

Introduction to scikit-Learn(sklearn)

This notebook demonstrates some of the most useful functions of the beautiful scikit library

What we are going to cover

0. An end-to-end scikit-learn workflow
1. Getting the data ready
2. Choose the right estimator/algorithm for our problems
3. Fit the models/algorithm and use it to make predictions on our data
4. Evaluating a model
5. Improve a model
6. Save and load a trained model
7. Putting it all together

```
In [12]: 1 import numpy as np
```

0. An end-to-end scikit-learn workflow

```
In [13]: 1 # 1. Getting the data ready
2 import pandas as pd
3 heart_disease = pd.read_csv('heart-disease.csv')
4 heart_disease
```

```
Out[13]:
```

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	63	1	3	145	233	1	0	150	0	2.3	0	0	1	1
1	37	1	2	130	250	0	1	187	0	3.5	0	0	2	1
2	41	0	1	130	204	0	0	172	0	1.4	2	0	2	1
3	56	1	1	120	236	0	1	178	0	0.8	2	0	2	1
4	57	0	0	120	354	0	1	163	1	0.6	2	0	2	1
...
298	57	0	0	140	241	0	1	123	1	0.2	1	0	3	0
299	45	1	3	110	264	0	1	132	0	1.2	1	0	3	0
300	68	1	0	144	193	1	1	141	0	3.4	1	2	3	0
301	57	1	0	130	131	0	1	115	1	1.2	1	1	3	0
302	57	0	1	130	236	0	0	174	0	0.0	1	1	2	0

303 rows × 14 columns

```
In [14]: 1 # create x (features matrix)
2 x = heart_disease.drop('target', axis = 1)
3
4 # create y (labels)
5 y = heart_disease['target']
6
```

```
In [15]: 1 # 2. choose the right model and hyperparameters
2 from sklearn.ensemble import RandomForestClassifier
3 clf = RandomForestClassifier()
4
5 # we'll keep the default hyperparameters
6 clf.get_params()
```

```
Out[15]: {'bootstrap': True,
'ccp_alpha': 0.0,
'class_weight': None,
'criterion': 'gini',
'max_depth': None,
'max_features': 'auto',
'max_leaf_nodes': None,
'max_samples': None,
'min_impurity_decrease': 0.0,
'min_samples_leaf': 1,
'min_samples_split': 2,
'min_weight_fraction_leaf': 0.0,
'n_estimators': 100,
'n_jobs': None,
'oob_score': False,
'random_state': None,
'verbose': 0,
'warm_start': False}
```

```
In [16]: 1 # Fit the model to the training data
2 from sklearn.model_selection import train_test_split
3
4 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0.2)
```

```
In [17]: 1 clf.fit(x_train, y_train)
```

```
Out[17]: RandomForestClassifier()
```

In [18]: 1 x_train

Out[18]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal
106	69	1	3	160	234	1	0	131	0	0.1	1	1	2
178	43	1	0	120	177	0	0	120	1	2.5	1	0	3
269	56	1	0	130	283	1	0	103	1	1.6	0	0	3
153	66	0	2	146	278	0	0	152	0	0.0	1	1	2
186	60	1	0	130	253	0	1	144	1	1.4	2	1	3
...
129	74	0	1	120	269	0	0	121	1	0.2	2	1	2
302	57	0	1	130	236	0	0	174	0	0.0	1	1	2
116	41	1	2	130	214	0	0	168	0	2.0	1	0	2
140	51	0	2	120	295	0	0	157	0	0.6	2	0	2
146	44	0	2	118	242	0	1	149	0	0.3	1	1	2

242 rows × 13 columns

```
In [19]: 1 # make a prediction  
2 y_label = clf.predict(np.array([0, 2, 3, 4]))
```

```
C:\Users\USER\anaconda3\lib\site-packages\sklearn\base.py:450: UserWarning: X does not have valid feature names, but RandomForestClassifier was fitted with feature names  
warnings.warn(
```

```

-----
ValueError                                Traceback (most recent call last)
~\AppData\Local\Temp\ipykernel_6928\3217224712.py in <module>
      1 # make a prediction
----> 2 y_label = clf.predict(np.array([0, 2, 3, 4]))

~\anaconda3\lib\site-packages\sklearn\ensemble\_forest.py in predict(self, X)
    806         The predicted classes.
    807         """
--> 808         proba = self.predict_proba(X)
    809
    810         if self.n_outputs_ == 1:

~\anaconda3\lib\site-packages\sklearn\ensemble\_forest.py in predict_proba(self, X)
    848         check_is_fitted(self)
    849         # Check data
--> 850         X = self._validate_X_predict(X)
    851
    852         # Assign chunk of trees to jobs

~\anaconda3\lib\site-packages\sklearn\ensemble\_forest.py in _validate_X_predict(self, X)
    577         Validate X whenever one tries to predict, apply, predict_proba."""
    578         check_is_fitted(self)
--> 579         X = self._validate_data(X, dtype=DTYPE, accept_sparse="csr", reset=False)
    580         if issparse(X) and (X.indices.dtype != np.intc or X.indptr.dtype != np.intc):
    581             raise ValueError("No support for np.int64 index based sparse matrices")

~\anaconda3\lib\site-packages\sklearn\base.py in _validate_data(self, X, y, reset, validate_separately,
**check_params)
    564             raise ValueError("Validation should be done on X, y or both.")
    565         elif not no_val_X and no_val_y:
--> 566             X = check_array(X, **check_params)
    567             out = X
    568         elif no_val_X and not no_val_y:

~\anaconda3\lib\site-packages\sklearn\utils\validation.py in check_array(array, accept_sparse, accept_large_sparse, dtype, order, copy, force_all_finite, ensure_2d, allow_nd, ensure_min_samples, ensure_min_features, estimator)
    767         # If input is 1D raise error
    768         if array.ndim == 1:
--> 769             raise ValueError(
    770                 "Expected 2D array, got 1D array instead:\nnarray={}\n"
    771                 "Reshape your data either using array.reshape(-1, 1) if "

```

ValueError: Expected 2D array, got 1D array instead:

array=[0. 2. 3. 4.].

Reshape your data either using `array.reshape(-1, 1)` if your data has a single feature or `array.reshape(1, -1)` if it contains a single sample.

In [20]:

```
1 x_test
```

Out[20]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal
192	54	1	0	120	188	0	1	113	0	1.4	1	1	3
63	41	1	1	135	203	0	1	132	0	0.0	1	0	1
20	59	1	0	135	234	0	1	161	0	0.5	1	0	3
273	58	1	0	100	234	0	1	156	0	0.1	2	1	3
164	38	1	2	138	175	0	1	173	0	0.0	2	4	2
...
237	60	1	0	140	293	0	0	170	0	1.2	1	2	3
179	57	1	0	150	276	0	0	112	1	0.6	1	1	1
194	60	1	2	140	185	0	0	155	0	3.0	1	0	2
241	59	0	0	174	249	0	1	143	1	0.0	1	0	2
124	39	0	2	94	199	0	1	179	0	0.0	2	0	2

61 rows × 13 columns

In [21]:

```
1 y_preds = clf.predict(x_test)
2 y_preds
```

Out[21]:

```
array([0, 1, 0, 1, 1, 0, 0, 0, 0, 1, 0, 1, 1, 0, 0, 0, 1, 1, 0, 1, 1, 1,
       0, 0, 1, 1, 1, 1, 0, 1, 0, 0, 1, 0, 0, 0, 1, 1, 1, 1, 0, 0, 0, 1,
       1, 0, 1, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 1, 0, 1], dtype=int64)
```

In [22]: 1 y_test

Out[22]: 192 0
63 1
20 1
273 0
164 1
..
237 0
179 0
194 0
241 0
124 1
Name: target, Length: 61, dtype: int64

In [23]: 1 # 4.Evaluate the model on the training data and test data
2 clf.score(x_train, y_train)

Out[23]: 1.0

In [24]: 1 clf.score(x_test, y_test)

Out[24]: 0.7540983606557377

In [25]: 1 from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
2 print(classification_report(y_test, y_preds))

	precision	recall	f1-score	support
0	0.84	0.73	0.78	37
1	0.66	0.79	0.72	24
accuracy			0.75	61
macro avg	0.75	0.76	0.75	61
weighted avg	0.77	0.75	0.76	61

In [26]: 1 confusion_matrix(y_test, y_preds)

Out[26]: array([[27, 10],
[5, 19]], dtype=int64)


```
In [27]: 1 accuracy_score(y_test, y_preds)
```

```
Out[27]: 0.7540983606557377
```

```
In [28]: 1 # improve a model
2 # try different amount of n_estimators
3 np.random.seed(42)
4 for i in range(10, 100, 10):
5     print(f"Trying model with {i} estimators...")
6     clf = RandomForestClassifier(n_estimators = i).fit(x_train, y_train)
7     print(f"Model accuracy on test set: {clf.score(x_test, y_test) * 100:.2f}%")
8     print("")
```

```
Trying model with 10 estimators...
Model accuracy on test set: 75.41%
```

```
Trying model with 20 estimators...
Model accuracy on test set: 77.05%
```

```
Trying model with 30 estimators...
Model accuracy on test set: 70.49%
```

```
Trying model with 40 estimators...
Model accuracy on test set: 72.13%
```

```
Trying model with 50 estimators...
Model accuracy on test set: 68.85%
```

```
Trying model with 60 estimators...
Model accuracy on test set: 72.13%
```

```
Trying model with 70 estimators...
Model accuracy on test set: 72.13%
```

```
Trying model with 80 estimators...
Model accuracy on test set: 73.77%
```

```
Trying model with 90 estimators...
Model accuracy on test set: 68.85%
```

```
In [29]: 1 # 6. save a model and load it
2 import pickle
3
4 pickle.dump(clf, open('random_forest_model_1.pk1', 'wb'))
```

```
In [30]: 1 loaded_model = pickle.load(open('random_forest_model_1.pk1', 'rb'))
2 loaded_model.score(x_test, y_test)
```

Out[30]: 0.6885245901639344

```
In [31]: 1 # to show the sklearn version we are using
2 import sklearn
3 sklearn.show_versions()
4
5 # to ignore this warning
6 import warnings
7 warnings.filterwarnings('ignore')
```

C:\Users\USER\anaconda3\lib\site-packages_distutils_hack__init__.py:33: UserWarning: Setuptools is replacing distutils.

warnings.warn("Setuptools is replacing distutils.")

System:

python: 3.9.13 (main, Aug 25 2022, 23:51:50) [MSC v.1916 64 bit (AMD64)]

executable: C:\Users\USER\anaconda3\python.exe

machine: Windows-10-10.0.19045-SP0

Python dependencies:

pip: 22.2.2

setuptools: 63.4.1

sklearn: 1.0.2

numpy: 1.21.5

scipy: 1.9.1

Cython: 0.29.32

pandas: 1.4.4

matplotlib: 3.5.2

joblib: 1.1.0

threadpoolctl: 2.2.0

Built with OpenMP: True

```
In [32]: 1 # Let's listify the contents
          2 what_were_covering = [
          3     '0. An end-to-end scikit-learn workflow',
          4     '1. Getting the data ready',
          5     '2. Choose the right estimator/algorithm for our problems',
          6     '3. Fit the models/algorithm and use it to make predictions on our data',
          7     '4. Evaluating a model',
          8     '5. Improve a model',
          9     '6. Save and load a trained model',
         10     '7. Putting it all together',
         11 ]
```

Second part of the lesson

```
In [33]: 1 what_were_covering
```

```
Out[33]: ['0. An end-to-end scikit-learn workflow',
          '1. Getting the data ready',
          '2. Choose the right estimator/algorithm for our problems',
          '3. Fit the models/algorithm and use it to make predictions on our data',
          '4. Evaluating a model',
          '5. Improve a model',
          '6. Save and load a trained model',
          '7. Putting it all together']
```

```
In [34]: 1 # standard imports
          2 import numpy as np
          3 import pandas as pd
          4 import matplotlib.pyplot as plt
          5 %matplotlib inline
```

1. Getting our data ready to be used with machine learning

Three main things we have to do:

1. Split the data into features and labels (usually 'x' & 'y')
2. Filling (also called imputing) or disregarding missing values
3. Converting non-numerical values to numerical values (also called feature encoding)

In [35]:

```
1 heart_disease.head()
```

Out[35]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	63	1	3	145	233	1	0	150	0	2.3	0	0	1	1
1	37	1	2	130	250	0	1	187	0	3.5	0	0	2	1
2	41	0	1	130	204	0	0	172	0	1.4	2	0	2	1
3	56	1	1	120	236	0	1	178	0	0.8	2	0	2	1
4	57	0	0	120	354	0	1	163	1	0.6	2	0	2	1

In [36]:

```
1 x = heart_disease.drop('target', axis = 1)
2 x.head()
```

Out[36]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal
0	63	1	3	145	233	1	0	150	0	2.3	0	0	1
1	37	1	2	130	250	0	1	187	0	3.5	0	0	2
2	41	0	1	130	204	0	0	172	0	1.4	2	0	2
3	56	1	1	120	236	0	1	178	0	0.8	2	0	2
4	57	0	0	120	354	0	1	163	1	0.6	2	0	2

In [37]:

```
1 y = heart_disease['target']
2 y.head()
```

Out[37]:

```
0    1
1    1
2    1
3    1
4    1
Name: target, dtype: int64
```

```
In [38]: 1 # split the data into training and test sets
2 from sklearn.model_selection import train_test_split
3 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0.2)
```

```
In [39]: 1 # checking the shape of our data
2 x_train.shape, x_test.shape, y_train.shape, y_test.shape
3
4 # here, test_size applied is = 0.2(20%) that is why roll is 242
```

Out[39]: ((242, 13), (61, 13), (242,), (61,))

```
In [40]: 1 # here is the full roll
2 x.shape
```

Out[40]: (303, 13)

```
In [41]: 1 # here is also the full roll
2 len(heart_disease)
```

Out[41]: 303

```
In [42]: 1 x.shape[0] * 0.8
```

Out[42]: 242.4

```
In [43]: 1 242 + 61
```

Out[43]: 303

make sure the data is all numerical

```
In [44]: 1 car_sales = pd.read_csv('car-sales-extended.csv')
         2 car_sales.head()
```

```
Out[44]:
```

	Make	Colour	Odometer (KM)	Doors	Price
0	Honda	White	35431	4	15323
1	BMW	Blue	192714	5	19943
2	Honda	White	84714	4	28343
3	Toyota	White	154365	4	13434
4	Nissan	Blue	181577	3	14043

```
In [45]: 1 car_sales['Doors'].value_counts()
```

```
Out[45]: 4    856
         5     79
         3     65
         Name: Doors, dtype: int64
```

```
In [46]: 1 len(car_sales)
```

```
Out[46]: 1000
```

```
In [47]: 1 car_sales.dtypes
```

```
Out[47]: Make          object
         Colour        object
         Odometer (KM)  int64
         Doors          int64
         Price          int64
         dtype: object
```

In [48]:

```
1 # split the data into 'x' and 'y'
2 x = car_sales.drop('Price', axis = 1)
3 y = car_sales['Price']
4
5 # split into training and test
6 from sklearn.model_selection import train_test_split
7 x_train, x_test, y_train, y_test = train_test_split(x,
8                                                    y,
9                                                    test_size = 0.2)
10
```

```
In [49]: 1 # build machine Learning
          2 from sklearn.ensemble import RandomForestRegressor
          3 model = RandomForestRegressor()
          4 model.fit(x_train, y_train)
          5 model.score(x_test, y_test)
```



```

-----
ValueError                                Traceback (most recent call last)
~\AppData\Local\Temp\ipykernel_6928\2751715000.py in <module>
      2 from sklearn.ensemble import RandomForestRegressor
      3 model = RandomForestRegressor()
----> 4 model.fit(x_train, y_train)
      5 model.score(x_test, y_test)

~\anaconda3\lib\site-packages\sklearn\ensemble\_forest.py in fit(self, X, y, sample_weight)
    325         if issparse(y):
    326             raise ValueError("sparse multilabel-indicator for y is not supported.")
--> 327         X, y = self._validate_data(
    328             X, y, multi_output=True, accept_sparse="csc", dtype=DTYPE
    329         )

~\anaconda3\lib\site-packages\sklearn\base.py in _validate_data(self, X, y, reset, validate_separately,
**check_params)
    579         y = check_array(y, **check_y_params)
    580     else:
--> 581         X, y = check_X_y(X, y, **check_params)
    582         out = X, y
    583

~\anaconda3\lib\site-packages\sklearn\utils\validation.py in check_X_y(X, y, accept_sparse, accept_large
_sparse, dtype, order, copy, force_all_finite, ensure_2d, allow_nd, multi_output, ensure_min_samples, en
sure_min_features, y_numeric, estimator)
    962         raise ValueError("y cannot be None")
    963
--> 964     X = check_array(
    965         X,
    966         accept_sparse=accept_sparse,

~\anaconda3\lib\site-packages\sklearn\utils\validation.py in check_array(array, accept_sparse, accept_la
rge_sparse, dtype, order, copy, force_all_finite, ensure_2d, allow_nd, ensure_min_samples, ensure_min_fe
atures, estimator)
    744         array = array.astype(dtype, casting="unsafe", copy=False)
    745     else:
--> 746         array = np.asarray(array, order=order, dtype=dtype)
    747     except ComplexWarning as complex_warning:
    748         raise ValueError(

~\anaconda3\lib\site-packages\pandas\core\generic.py in __array__(self, dtype)
    2062

```

```

2063     def __array__(self, dtype: npt.DTypeLike | None = None) -> np.ndarray:
-> 2064         return np.asarray(self._values, dtype=dtype)
2065
2066     def __array_wrap__(

```

ValueError: could not convert string to float: 'Toyota'

In [50]:

```

1 # to know what x is before convrsion
2 x.head()

```

Out[50]:

	Make	Colour	Odometer (KM)	Doors
0	Honda	White	35431	4
1	BMW	Blue	192714	5
2	Honda	White	84714	4
3	Toyota	White	154365	4
4	Nissan	Blue	181577	3

```
In [51]: 1 # to convert the strings/objects to integer(numbers)(make and colour) from the data set to solve the
2 from sklearn.preprocessing import OneHotEncoder
3 from sklearn.compose import ColumnTransformer
4
5 # find the category features to convert to numbers
6 categorical_features = ['Make', 'Colour', 'Doors']
7 one_hot = OneHotEncoder()
8 transformer = ColumnTransformer([('one_hot',
9                                   one_hot,
10                                  categorical_features)],
11                                remainder = 'passthrough')
12
13 transformed_x = transformer.fit_transform(x)
14 transformed_x
```

```
Out[51]: array([[0.00000e+00, 1.00000e+00, 0.00000e+00, ..., 1.00000e+00,
0.00000e+00, 3.54310e+04],
[1.00000e+00, 0.00000e+00, 0.00000e+00, ..., 0.00000e+00,
1.00000e+00, 1.92714e+05],
[0.00000e+00, 1.00000e+00, 0.00000e+00, ..., 1.00000e+00,
0.00000e+00, 8.47140e+04],
...,
[0.00000e+00, 0.00000e+00, 1.00000e+00, ..., 1.00000e+00,
0.00000e+00, 6.66040e+04],
[0.00000e+00, 1.00000e+00, 0.00000e+00, ..., 1.00000e+00,
0.00000e+00, 2.15883e+05],
[0.00000e+00, 0.00000e+00, 0.00000e+00, ..., 1.00000e+00,
0.00000e+00, 2.48360e+05]])
```

```
In [52]: 1 # putting the above data in a data frame
          2
          3 pd.DataFrame(transformed_x).head()
```

```
Out[52]:
```

	0	1	2	3	4	5	6	7	8	9	10	11	12
0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	1.0	0.0	35431.0
1	1.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	1.0	192714.0
2	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	1.0	0.0	84714.0
3	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	1.0	0.0	1.0	0.0	154365.0
4	0.0	0.0	1.0	0.0	0.0	1.0	0.0	0.0	0.0	1.0	0.0	0.0	181577.0

```
In [53]: 1 # another way to transform our column to integers(numbers)
          2
          3 dummies = pd.get_dummies(car_sales[['Make', 'Colour', 'Doors']])
          4 dummies.head()
```

```
Out[53]:
```

	Doors	Make_BMW	Make_Honda	Make_Nissan	Make_Toyota	Colour_Black	Colour_Blue	Colour_Green	Colour_Red	Colour_White
0	4	0	1	0	0	0	0	0	0	1
1	5	1	0	0	0	0	1	0	0	0
2	4	0	1	0	0	0	0	0	0	1
3	4	0	0	0	1	0	0	0	0	1
4	3	0	0	1	0	0	1	0	0	0

```
In [54]: 1 # since our data has turn to zero and ones
          2 # Let's refill the model
          3
          4 np.random.seed(42)
          5 x_train, x_test, y_train, y_test = train_test_split(transformed_x,
          6                                                         y,
          7                                                         test_size=0.2)
          8 model.fit(x_train, y_train)
```

```
Out[54]: RandomForestRegressor()
```

```
In [55]: 1 model.score(x_test, y_test)
```

```
Out[55]: 0.3235867221569877
```

1.2 what if there were missing values?

1. Fill them with some value (also known as imputation).
2. Remove the samples with missing data altogether.

```
In [56]: 1 # import car sales missind data
2 car_sales_missing = pd.read_csv('car-sales-extended-missing-data.csv')
3 car_sales_missing.head()
```

```
Out[56]:
```

	Make	Colour	Odometer (KM)	Doors	Price
0	Honda	White	35431.0	4.0	15323.0
1	BMW	Blue	192714.0	5.0	19943.0
2	Honda	White	84714.0	4.0	28343.0
3	Toyota	White	154365.0	4.0	13434.0
4	Nissan	Blue	181577.0	3.0	14043.0

```
In [57]: 1 # to check missing data
2 car_sales_missing.isna().sum()
```

```
Out[57]: Make          49
Colour          50
Odometer (KM)   50
Doors           50
Price           50
dtype: int64
```

```
In [58]: 1 # create x & y
2 x = car_sales_missing.drop('Price', axis = 1)
3 y = car_sales_missing['Price']
```

```

In [59]: 1 # Let's try and convert our data to numbers
          2 from sklearn.preprocessing import OneHotEncoder
          3 from sklearn.compose import ColumnTransformer
          4
          5 # find the category features to convert to numbers
          6 categorical_features = ['Make', 'Colour', 'Doors']
          7 one_hot = OneHotEncoder()
          8 transformer = ColumnTransformer([('one_hot',
          9                                one_hot,
          10                                categorical_features)],
          11                                remainder = 'passthrough')
          12
          13 transformed_x = transformer.fit_transform(x)
          14 transformed_x

```

Out[59]: <1000x16 sparse matrix of type '<class 'numpy.float64'>'
with 4000 stored elements in Compressed Sparse Row format>

```

In [60]: 1 car_sales_missing

```

Out[60]:

	Make	Colour	Odometer (KM)	Doors	Price
0	Honda	White	35431.0	4.0	15323.0
1	BMW	Blue	192714.0	5.0	19943.0
2	Honda	White	84714.0	4.0	28343.0
3	Toyota	White	154365.0	4.0	13434.0
4	Nissan	Blue	181577.0	3.0	14043.0
...
995	Toyota	Black	35820.0	4.0	32042.0
996	NaN	White	155144.0	3.0	5716.0
997	Nissan	Blue	66604.0	4.0	31570.0
998	Honda	White	215883.0	4.0	4001.0
999	Toyota	Blue	248360.0	4.0	12732.0

1000 rows × 5 columns

```
In [61]: 1 # above worked directly without filling up with pandas
```

```
In [62]: 1 car_sales_missing['Doors'].value_counts()
```

```
Out[62]: 4.0    811  
5.0     75  
3.0     64  
Name: Doors, dtype: int64
```

option 1: fill missing data with pandas

```
In [63]: 1 # fill the 'Make' column  
2 car_sales_missing['Make'].fillna('missing', inplace=True)  
3  
4 # fill the 'Colour' column  
5 car_sales_missing['Colour'].fillna('missing', inplace = True)  
6  
7 # fill the 'Odometer (KM)' column  
8 car_sales_missing['Odometer (KM)'].fillna(car_sales_missing['Odometer (KM)'].mean(), inplace = True)  
9  
10 # fill the 'Doors' column  
11 car_sales_missing['Doors'].fillna(4, inplace = True)
```

```
In [64]: 1 # check our dataframe again  
2 car_sales_missing.isna().sum()
```

```
Out[64]: Make          0  
Colour          0  
Odometer (KM)    0  
Doors           0  
Price          50  
dtype: int64
```

```
In [65]: 1 # remove rows with missing price values  
2 car_sales_missing.dropna(inplace = True)  
3
```

```
In [66]: 1 car_sales_missing.isna().sum()
```

```
Out[66]: Make          0  
Colour         0  
Odometer (KM)   0  
Doors          0  
Price          0  
dtype: int64
```

```
In [67]: 1 len(car_sales_missing)
```

```
Out[67]: 950
```

```
In [68]: 1 x = car_sales_missing.drop('Price', axis = 1)  
2 y = car_sales_missing['Price']
```



```
In [69]: 1 # Let's bring back our code now
2
3 # Let's try and convert our data to numbers
4 from sklearn.preprocessing import OneHotEncoder
5 from sklearn.compose import ColumnTransformer
6
7 # find the category features to convert to numbers
8 categorical_features = ['Make', 'Colour', 'Doors']
9 one_hot = OneHotEncoder()
10 transformer = ColumnTransformer([('one_hot',
11                                 one_hot,
12                                 categorical_features)],
13                                remainder = 'passthrough')
14
15 transformed_x = transformer.fit_transform(car_sales_missing)
16 transformed_x
```

```
Out[69]: array([[0.00000e+00, 1.00000e+00, 0.00000e+00, ..., 0.00000e+00,
3.54310e+04, 1.53230e+04],
[1.00000e+00, 0.00000e+00, 0.00000e+00, ..., 1.00000e+00,
1.92714e+05, 1.99430e+04],
[0.00000e+00, 1.00000e+00, 0.00000e+00, ..., 0.00000e+00,
8.47140e+04, 2.83430e+04],
...,
[0.00000e+00, 0.00000e+00, 1.00000e+00, ..., 0.00000e+00,
6.66040e+04, 3.15700e+04],
[0.00000e+00, 1.00000e+00, 0.00000e+00, ..., 0.00000e+00,
2.15883e+05, 4.00100e+03],
[0.00000e+00, 0.00000e+00, 0.00000e+00, ..., 0.00000e+00,
2.48360e+05, 1.27320e+04]])
```

```
In [70]: 1 # putting above data into a dataframe
          2 pd.DataFrame(transformed_x).head()
```

```
Out[70]:
```

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	1.0	0.0	35431.0	15323.0
1	1.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	192714.0	19943.0
2	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	1.0	0.0	84714.0	28343.0
3	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	1.0	0.0	154365.0	13434.0
4	0.0	0.0	1.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	181577.0	14043.0

Option 2: fill missing value with scikit-learn

```
In [71]: 1 car_sales_missing = pd.read_csv('car-sales-extended-missing-data.csv')
          2 car_sales_missing.head()
```

```
Out[71]:
```

	Make	Colour	Odometer (KM)	Doors	Price
0	Honda	White	35431.0	4.0	15323.0
1	BMW	Blue	192714.0	5.0	19943.0
2	Honda	White	84714.0	4.0	28343.0
3	Toyota	White	154365.0	4.0	13434.0
4	Nissan	Blue	181577.0	3.0	14043.0

```
In [72]: 1 # to check for missing data
          2 car_sales_missing.isna().sum()
```

```
Out[72]: Make          49
          Colour       50
          Odometer (KM) 50
          Doors        50
          Price        50
          dtype: int64
```

```
In [73]: 1 car_sales_missing.dropna(subset = ['Price'], inplace = True)
        2 car_sales_missing
```

```
Out[73]:
```

	Make	Colour	Odometer (KM)	Doors	Price
0	Honda	White	35431.0	4.0	15323.0
1	BMW	Blue	192714.0	5.0	19943.0
2	Honda	White	84714.0	4.0	28343.0
3	Toyota	White	154365.0	4.0	13434.0
4	Nissan	Blue	181577.0	3.0	14043.0
...
995	Toyota	Black	35820.0	4.0	32042.0
996	NaN	White	155144.0	3.0	5716.0
997	Nissan	Blue	66604.0	4.0	31570.0
998	Honda	White	215883.0	4.0	4001.0
999	Toyota	Blue	248360.0	4.0	12732.0

950 rows × 5 columns

```
In [74]: 1 car_sales_missing.isna().sum()
```

```
Out[74]: Make          47
         Colour        46
         Odometer (KM)  48
         Doors         47
         Price          0
         dtype: int64
```

```
In [75]: 1 # split into x and y
        2 x = car_sales_missing.drop('Price', axis = 1)
        3 y = car_sales_missing['Price']
```

```
In [76]: 1 x.isna().sum()
```

```
Out[76]: Make                47  
Colour                46  
Odometer (KM)        48  
Doors                 47  
dtype: int64
```

```
In [77]: 1 # to fix missing data/values with sklearn  
2  
3 from sklearn.impute import SimpleImputer  
4 from sklearn.compose import ColumnTransformer  
5  
6 # fill categorical values 'missing' & numerical values with 'mean'  
7 cat_imputer = SimpleImputer(strategy = 'constant', fill_value = 'missing')  
8 door_imputer = SimpleImputer(strategy = 'constant', fill_value = 4)  
9 num_imputer = SimpleImputer(strategy = 'mean')  
10  
11 # Define columns  
12 cat_features = ['Make', 'Colour']  
13 door_feature = ['Doors']  
14 num_features = ['Odometer (KM)']  
15  
16 # create an imputer (something that fills missing data)  
17 imputer = ColumnTransformer([  
18     ('cat_imputer', cat_imputer, cat_features),  
19     ('door_imputer', door_imputer, door_feature),  
20     ('num_imputer', num_imputer, num_features)  
21 ])  
22  
23 # transform the data  
24 filled_x = imputer.fit_transform(x)  
25 filled_x  
26
```

```
Out[77]: array([[ 'Honda', 'White', 4.0, 35431.0],  
                [ 'BMW', 'Blue', 5.0, 192714.0],  
                [ 'Honda', 'White', 4.0, 84714.0],  
                ...,  
                [ 'Nissan', 'Blue', 4.0, 66604.0],  
                [ 'Honda', 'White', 4.0, 215883.0],  
                [ 'Toyota', 'Blue', 4.0, 248360.0]], dtype=object)
```

```
In [78]: 1 car_sales_filled = pd.DataFrame(filled_x,  
2                                     columns = ['Make', 'Colour', 'Doors', 'Odometer (KM)'])  
3 car_sales_filled.head()
```

```
Out[78]:
```

	Make	Colour	Doors	Odometer (KM)
0	Honda	White	4.0	35431.0
1	BMW	Blue	5.0	192714.0
2	Honda	White	4.0	84714.0
3	Toyota	White	4.0	154365.0
4	Nissan	Blue	3.0	181577.0

```
In [79]: 1 car_sales_filled.isna().sum()
```

```
Out[79]: Make          0  
Colour          0  
Doors           0  
Odometer (KM)   0  
dtype: int64
```

```
In [80]: 1 car_sales_filled.head()
```

```
Out[80]:
```

	Make	Colour	Doors	Odometer (KM)
0	Honda	White	4.0	35431.0
1	BMW	Blue	5.0	192714.0
2	Honda	White	4.0	84714.0
3	Toyota	White	4.0	154365.0
4	Nissan	Blue	3.0	181577.0

```
In [81]: 1 # Let's try and convert our data to numbers
2 from sklearn.preprocessing import OneHotEncoder
3 from sklearn.compose import ColumnTransformer
4
5 # find the category features to convert to numbers
6 categorical_features = ['Make', 'Colour', 'Doors']
7 one_hot = OneHotEncoder()
8 transformer = ColumnTransformer([('one_hot',
9                                 one_hot,
10                                categorical_features)],
11                                remainder = 'passthrough')
12
13 transformed_x = transformer.fit_transform(car_sales_filled)
14 transformed_x
```

```
Out[81]: <950x15 sparse matrix of type '<class 'numpy.float64'>'
         with 3800 stored elements in Compressed Sparse Row format>
```

```
In [82]: 1 # Now we've got our data as numbers and filled (no missing values)
2 # Let's fit a model
3 np.random.seed(42)
4 from sklearn.ensemble import RandomForestRegressor
5 from sklearn.model_selection import train_test_split
6
7 x_train, x_test, y_train, y_test = train_test_split(transformed_x,
8                                                    y,
9                                                    test_size = 0.2)
10 model = RandomForestRegressor()
11 model.fit(x_train, y_train)
12 model.score(x_test, y_test)
13
```

```
Out[82]: 0.21990196728583944
```

```
In [83]: 1 # the above performed worse cos it has 950 samples
2 len(car_sales_filled), len(car_sales)
```

```
Out[83]: (950, 1000)
```

2. Choosing the right estimator/algorithm for our problem

scikit-learn uses estimator as another term for machine learning model or algorithm.

- Classification - predicting whether a sample is one thing or another.
- Regression - predicting a number

step 1 - check the scikit-learn machine learning mapp... https://scikit-learn.org/stable/tutorial/machine_learning_map/index.html
(https://scikit-learn.org/stable/tutorial/machine_learning_map/index.html).

2.1 picking a machine learning model for a regression problem

```
In [84]: 1 # import Boston housing dataset
          2 from sklearn.datasets import load_boston
          3 boston = load_boston()
          4 boston;
```

```
In [85]: 1 boston_df = pd.DataFrame(boston['data'], columns = boston['feature_names'])
          2 boston_df['target'] = pd.Series(boston['target'])
          3 boston_df.head()
```

```
Out[85]:
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	target
0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1.0	296.0	15.3	396.90	4.98	24.0
1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2.0	242.0	17.8	396.90	9.14	21.6
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2.0	242.0	17.8	392.83	4.03	34.7
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3.0	222.0	18.7	394.63	2.94	33.4
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3.0	222.0	18.7	396.90	5.33	36.2

```
In [86]: 1 # how many samples?
          2 len(boston_df)
```

```
Out[86]: 506
```

```
In [87]: 1 # Let try the Ridge Regression model
2 from sklearn.linear_model import Ridge
3
4 # setup random seed
5 np.random.seed(42)
6
7 # create the data
8 x = boston_df.drop('target', axis = 1)
9 y = boston_df['target']
10
11 # split into train and test sets
12
13 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0.2)
14
15 # instantiate Ridge model
16 model = Ridge()
17 model.fit(x_train, y_train)
18
19 # checck the score of the Ridge
20 model.score(x_test, y_test)
```

Out[87]: 0.6662221670168518

How do we improve this score?

What if Ridge wasn't working

let's refer back to the map... https://scikit-learn.org/stable/tutorial/machine_learning_map/index.html (https://scikit-learn.org/stable/tutorial/machine_learning_map/index.html).


```
In [88]: 1 # Let's try the Random Forest Regressor
2 from sklearn.ensemble import RandomForestRegressor
3
4 # set random seed
5 np.random.seed(42)
6
7 # create the data
8 x = boston_df.drop('target', axis = 1)
9 y = boston_df['target']
10
11 # split the data
12 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)
13
14 # instantiate random forest regressor
15 rf = RandomForestRegressor()
16 rf.fit(x_train, y_train)
17
18 # evaluate the random forest regressor
19 rf.score(x_test, y_test)
20
```

Out[88]: 0.8654448653350507

```
In [89]: 1 # ccheck the Ridge model again
2 model.score(x_test, y_test)
3
```

Out[89]: 0.6662221670168518

Choosing an estimator for a classification problem

let's go tp the maap...https://scikit-learn.org/stable/tutorial/machine_learning_map/index.html (https://scikit-learn.org/stable/tutorial/machine_learning_map/index.html)

```
In [90]: 1 heart_disease = pd.read_csv('heart-disease.csv')
        2 heart_disease.head()
```

```
Out[90]:
```

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	63	1	3	145	233	1	0	150	0	2.3	0	0	1	1
1	37	1	2	130	250	0	1	187	0	3.5	0	0	2	1
2	41	0	1	130	204	0	0	172	0	1.4	2	0	2	1
3	56	1	1	120	236	0	1	178	0	0.8	2	0	2	1
4	57	0	0	120	354	0	1	163	1	0.6	2	0	2	1

```
In [91]: 1 len(heart_disease)
```

```
Out[91]: 303
```

Consulting the map and it says to try linearSVC.

```
In [92]: 1 # import the LinearSVC estimator class
2 from sklearn.svm import LinearSVC
3
4 # set up random seed
5 np.random.seed(42)
6
7 # make the data
8 x = heart_disease.drop('target', axis = 1)
9 y = heart_disease['target']
10
11 # split the data
12 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)
13
14 # instantiate LinearSVC
15 ln = LinearSVC()
16 ln.fit(x_train, y_train)
17
18 #Evaluate the LinearSVC
19 ln.score(x_test, y_test)
20
```

Out[92]: 0.8688524590163934

```
In [93]: 1 heart_disease['target'].value_counts()
```

Out[93]: 1 165
0 138
Name: target, dtype: int64

```

In [94]: 1 # import the RandomForestClassifier estimator class
          2 from sklearn.ensemble import RandomForestClassifier
          3
          4 # set up random seed
          5 np.random.seed(42)
          6
          7 # make the data
          8 x = heart_disease.drop('target', axis = 1)
          9 y = heart_disease['target']
         10
         11 # split the data
         12 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)
         13
         14 # instantiate RandomForestClassifier
         15 rfc = RandomForestClassifier()
         16 rfc.fit(x_train, y_train)
         17
         18 #Evaluate the RandomForestClassifier
         19 rfc.score(x_test, y_test)
         20

```

Out[94]: 0.8524590163934426

Tidbit:

1. if you have structured(tables) data, use ensemble methods
2. if you have unstructured(images, audios etc) data, use deep learning or transfer learning

```

In [95]: 1 heart_disease.head() # this is a structured data

```

```

Out[95]:
   age  sex  cp  trestbps  chol  fbs  restecg  thalach  exang  oldpeak  slope  ca  thal  target
0   63    1   3    145    233    1         0     150      0        2.3     0   0    1         1
1   37    1   2    130    250    0         1     187      0        3.5     0   0    2         1
2   41    0   1    130    204    0         0     172      0        1.4     2   0    2         1
3   56    1   1    120    236    0         1     178      0        0.8     2   0    2         1
4   57    0   0    120    354    0         1     163      1        0.6     2   0    2         1

```

3. Fit the model/algorithm on our data and use it to make predictions

3.1 Fitting the model to the data

Different names for:

- x = features, features variables, data
- y = labels, targets, target variables

In [96]:

```
1  # import the RandomForestClassifier estimator class
2  from sklearn.ensemble import RandomForestClassifier
3
4  # set up random seed
5  np.random.seed(42)
6
7  # make the data
8  x = heart_disease.drop('target', axis = 1)
9  y = heart_disease['target']
10
11 # split the data
12 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)
13
14 # instantiate RandomForestClassifier
15 rfc = RandomForestClassifier()
16
17 # fit the model to the data (training the machine learning model)
18 rfc.fit(x_train, y_train)
19
20 #Evaluate the RandomForestClassifier (use the patterns the model has learned)
21 rfc.score(x_test, y_test)
```

Out[96]: 0.8524590163934426

In [97]: 1 x.head()

Out[97]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal
0	63	1	3	145	233	1	0	150	0	2.3	0	0	1
1	37	1	2	130	250	0	1	187	0	3.5	0	0	2
2	41	0	1	130	204	0	0	172	0	1.4	2	0	2
3	56	1	1	120	236	0	1	178	0	0.8	2	0	2
4	57	0	0	120	354	0	1	163	1	0.6	2	0	2

In [98]: 1 y.head()

Out[98]:

0	1
1	1
2	1
3	1
4	1

Name: target, dtype: int64

In [99]: 1 y.tail()

Out[99]:

298	0
299	0
300	0
301	0
302	0

Name: target, dtype: int64

3.2 make predictions using a machine learning model

Two ways to make predictions

1. predict()
2. predict_proba()

In [100]:

```
1 # use a trained model to make predictions
2
3 rfc.predict(np.array([1, 7, 8, 3, 4])) # this doesn't work..
```

```

-----
ValueError                                Traceback (most recent call last)
~\AppData\Local\Temp\ipykernel_6928\4195040950.py in <module>
      1 # use a trained model to make predictions
      2
----> 3 rfc.predict(np.array([1, 7, 8, 3, 4])) # this doesn't work..

~\anaconda3\lib\site-packages\sklearn\ensemble\_forest.py in predict(self, X)
    806         The predicted classes.
    807         """
--> 808         proba = self.predict_proba(X)
    809
    810         if self.n_outputs_ == 1:

~\anaconda3\lib\site-packages\sklearn\ensemble\_forest.py in predict_proba(self, X)
    848         check_is_fitted(self)
    849         # Check data
--> 850         X = self._validate_X_predict(X)
    851
    852         # Assign chunk of trees to jobs

~\anaconda3\lib\site-packages\sklearn\ensemble\_forest.py in _validate_X_predict(self, X)
    577         Validate X whenever one tries to predict, apply, predict_proba."""
    578         check_is_fitted(self)
--> 579         X = self._validate_data(X, dtype=DTYPE, accept_sparse="csr", reset=False)
    580         if issparse(X) and (X.indices.dtype != np.intc or X.indptr.dtype != np.intc):
    581             raise ValueError("No support for np.int64 index based sparse matrices")

~\anaconda3\lib\site-packages\sklearn\base.py in _validate_data(self, X, y, reset, validate_separately,
**check_params)
    564             raise ValueError("Validation should be done on X, y or both.")
    565             elif not no_val_X and no_val_y:
--> 566             X = check_array(X, **check_params)
    567             out = X
    568             elif no_val_X and not no_val_y:

~\anaconda3\lib\site-packages\sklearn\utils\validation.py in check_array(array, accept_sparse, accept_large_sparse, dtype, order, copy, force_all_finite, ensure_2d, allow_nd, ensure_min_samples, ensure_min_features, estimator)
    767         # If input is 1D raise error
    768         if array.ndim == 1:
--> 769             raise ValueError(
    770                 "Expected 2D array, got 1D array instead:\nnarray={}\n"

```


771

"Reshape your data either using array.reshape(-1, 1) if "**ValueError:** Expected 2D array, got 1D array instead:

array=[1. 7. 8. 3. 4.].

Reshape your data either using array.reshape(-1, 1) if your data has a single feature or array.reshape(1, -1) if it contains a single sample.

In [101]: 1 x_test.head()

Out[101]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal
179	57	1	0	150	276	0	0	112	1	0.6	1	1	1
228	59	1	3	170	288	0	0	159	0	0.2	1	0	3
111	57	1	2	150	126	1	1	173	0	0.2	2	1	3
246	56	0	0	134	409	0	0	150	1	1.9	1	2	3
60	71	0	2	110	265	1	0	130	0	0.0	2	1	2

In [102]: 1 rfc.predict(x_test)

Out[102]: array([0, 1, 1, 0, 1, 1, 1, 0, 0, 1, 1, 0, 1, 0, 1, 1, 1, 0, 0, 0, 0, 0,
1, 1, 1, 1, 1, 1, 0, 1, 0, 0, 0, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1,
1, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0], dtype=int64)

In [103]: 1 y_test

Out[103]: 179 0
228 0
111 1
246 0
60 1
..
249 0
104 1
300 0
193 0
184 0
Name: target, Length: 61, dtype: int64

```
In [104]: 1 # we can put y_test in form of x-test array
          2
          3 np.array(y_test)
```

```
Out[104]: array([0, 0, 1, 0, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 1, 1, 1, 0, 0, 0, 1, 0,
                0, 1, 1, 0, 0, 1, 0, 1, 1, 0, 0, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1,
                1, 0, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0], dtype=int64)
```

```
In [105]: 1 # compare predictions to truth labels to evaluate the model
          2 y_preds = rfc.predict(x_test)
          3 np.mean(y_preds == y_test)
```

```
Out[105]: 0.8524590163934426
```

```
In [106]: 1 rfc.score(x_test, y_test)
```

```
Out[106]: 0.8524590163934426
```

```
In [107]: 1 # another way of doing the above
          2 from sklearn.metrics import accuracy_score
          3 accuracy_score(y_test, y_preds)
```

```
Out[107]: 0.8524590163934426
```

Make predictions with predict_proba()

```
In [108]: 1 # Predicts_proba() returns probabilities of a classification label
          2 rfc.predict_proba(x_test[: 5])
          3
```

```
Out[108]: array([[0.89, 0.11],
                [0.49, 0.51],
                [0.43, 0.57],
                [0.84, 0.16],
                [0.18, 0.82]])
```

```
In [109]: 1 # Let's predict() on the same data...
          2 rfc.predict(x_test [: 5])
```

```
Out[109]: array([0, 1, 1, 0, 1], dtype=int64)
```

```
In [110]: 1 x_test[: 5]
```

```
Out[110]:
```

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal
179	57	1	0	150	276	0	0	112	1	0.6	1	1	1
228	59	1	3	170	288	0	0	159	0	0.2	1	0	3
111	57	1	2	150	126	1	1	173	0	0.2	2	1	3
246	56	0	0	134	409	0	0	150	1	1.9	1	2	3
60	71	0	2	110	265	1	0	130	0	0.0	2	1	2

```
In [111]: 1 heart_disease['target'].value_counts()
```

```
Out[111]: 1    165  
          0    138  
          Name: target, dtype: int64
```

making prediction with Regression Model

predict() can also be used for regression models

In [112]:

1 boston_df

Out[112]:

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	target
0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1.0	296.0	15.3	396.90	4.98	24.0
1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2.0	242.0	17.8	396.90	9.14	21.6
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2.0	242.0	17.8	392.83	4.03	34.7
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3.0	222.0	18.7	394.63	2.94	33.4
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3.0	222.0	18.7	396.90	5.33	36.2
...
501	0.06263	0.0	11.93	0.0	0.573	6.593	69.1	2.4786	1.0	273.0	21.0	391.99	9.67	22.4
502	0.04527	0.0	11.93	0.0	0.573	6.120	76.7	2.2875	1.0	273.0	21.0	396.90	9.08	20.6
503	0.06076	0.0	11.93	0.0	0.573	6.976	91.0	2.1675	1.0	273.0	21.0	396.90	5.64	23.9
504	0.10959	0.0	11.93	0.0	0.573	6.794	89.3	2.3889	1.0	273.0	21.0	393.45	6.48	22.0
505	0.04741	0.0	11.93	0.0	0.573	6.030	80.8	2.5050	1.0	273.0	21.0	396.90	7.88	11.9

506 rows × 14 columns

```
In [113]: 1 from sklearn.ensemble import RandomForestRegressor
2
3 np.random.seed(42)
4
5 # create the data
6 x = boston_df.drop('target', axis = 1)
7 y = boston_df['target']
8
9 # split into training and test sets
10 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)
11
12 # instantiate and fit the model
13 model = RandomForestRegressor().fit(x_train, y_train)
14
15 # make predictions
16
17 y_preds = model.predict(x_test)
```

```
In [114]: 1 y_preds[:10]
```

```
Out[114]: array([23.081, 30.574, 16.759, 23.46 , 16.893, 21.644, 19.113, 15.334,
                21.14 , 20.639])
```

```
In [115]: 1 np.array(y_test[:10])
```

```
Out[115]: array([23.6, 32.4, 13.6, 22.8, 16.1, 20. , 17.8, 14. , 19.6, 16.8])
```

```
In [116]: 1 # compare the predictions to the truth
2 from sklearn.metrics import mean_absolute_error
3 mean_absolute_error(y_test, y_preds)
```

```
Out[116]: 2.136382352941176
```

4. Evaluating a machine learning model(score)

Three ways to evaluate scikit-learn models/estimators:

1. Estimator score method
2. The scoring parameter
3. Problem specific metric functions

4.1 Evaluating a model with the 'score method'

In [117]:

```
1 heart_disease.head()
```

Out[117]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	63	1	3	145	233	1	0	150	0	2.3	0	0	1	1
1	37	1	2	130	250	0	1	187	0	3.5	0	0	2	1
2	41	0	1	130	204	0	0	172	0	1.4	2	0	2	1
3	56	1	1	120	236	0	1	178	0	0.8	2	0	2	1
4	57	0	0	120	354	0	1	163	1	0.6	2	0	2	1

In [118]:

```
1 from sklearn.ensemble import RandomForestClassifier
2
3 np.random.seed(42)
4
5 x = heart_disease.drop('target', axis = 1)
6 y = heart_disease['target']
7
8
9 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)
10
11 rfc = RandomForestClassifier()
12
13 rfc.fit(x_train, y_train)
```

Out[118]: RandomForestClassifier()

In [119]:

```
1 rfc.score(x_train, y_train)
```

Out[119]: 1.0

In [120]:

```
1 rfc.score(x_test, y_test)
```

Out[120]: 0.8524590163934426

let's do the same but for regression...

```
In [121]: 1 from sklearn.ensemble import RandomForestRegressor
          2
          3 np.random.seed(42)
          4
          5 # create the data
          6 x = boston_df.drop('target', axis = 1)
          7 y = boston_df['target']
          8
          9 # split into training and test sets
         10 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)
         11
         12 # instantiate and fit the model
         13 model = RandomForestRegressor().fit(x_train, y_train)
         14
```

```
In [122]: 1 model.score(x_test, y_test)
```

```
Out[122]: 0.8654448653350507
```

```
In [123]: 1 model.score(x_train, y_train)
```

```
Out[123]: 0.9763520974033731
```

4.2 Evaluating a model using scoring parameters (cross validation)

```
In [124]: 1 from sklearn.model_selection import cross_val_score
2
3 from sklearn.ensemble import RandomForestClassifier
4
5 np.random.seed(42)
6
7 x = heart_disease.drop('target', axis = 1)
8 y = heart_disease['target']
9
10
11 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)
12
13 rfc = RandomForestClassifier()
14
15 rfc.fit(x_train, y_train);
16
```

```
In [125]: 1 rfc.score(x_train, y_train)
```

```
Out[125]: 1.0
```

```
In [126]: 1 rfc.score(x_test, y_test)
```

```
Out[126]: 0.8524590163934426
```

```
In [127]: 1 cross_val_score(rfc, x, y)
```

```
Out[127]: array([0.81967213, 0.86885246, 0.81967213, 0.78333333, 0.76666667])
```

```
In [128]: 1 cross_val_score(rfc, x, y, cv = 10)
```

```
Out[128]: array([0.90322581, 0.80645161, 0.87096774, 0.9          , 0.86666667,
                  0.8          , 0.73333333, 0.86666667, 0.73333333, 0.8          ])
```



```
In [129]: 1 np.random.seed(42)
          2
          3 # single training and test split score
          4 rfc_single_score = rfc.score(x_test, y_test)
          5
          6 # take the mean of 5-fold cross validation score
          7 rfc_cross_val_score = np.mean(cross_val_score(rfc, x, y))
          8
          9 # compare the two
         10 rfc_single_score, rfc_cross_val_score
```

Out[129]: (0.8524590163934426, 0.8248087431693989)

```
In [130]: 1 # default scoring parameter of classifier = mean accuracy
          2 rfc.score(x, y)
```

Out[130]: 0.9702970297029703

```
In [131]: 1 # scoring parameter set to None by default
          2 cross_val_score(rfc, x, y, scoring = None)
```

Out[131]: array([0.78688525, 0.86885246, 0.80327869, 0.78333333, 0.76666667])

```
In [132]: 1 rfc.score(x_train, y_train)
```

Out[132]: 1.0

4.2.1 Classification model evaluation metrics

1. Accuracy
2. Area under ROC curve
3. confusion matrix
4. classification report

Accuracy

In [133]: 1 heart_disease.head()

Out[133]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	63	1	3	145	233	1	0	150	0	2.3	0	0	1	1
1	37	1	2	130	250	0	1	187	0	3.5	0	0	2	1
2	41	0	1	130	204	0	0	172	0	1.4	2	0	2	1
3	56	1	1	120	236	0	1	178	0	0.8	2	0	2	1
4	57	0	0	120	354	0	1	163	1	0.6	2	0	2	1

In [134]:

```

1 from sklearn.model_selection import cross_val_score
2 from sklearn.ensemble import RandomForestClassifier
3
4 np.random.seed(42)
5
6 x = heart_disease.drop('target', axis = 1)
7 y = heart_disease['target']
8
9 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)
10 clf = RandomForestClassifier()
11 clf.fit(x_train, y_train)
12 cross_val_score = cross_val_score(clf, x, y, cv = 5)

```

In [135]: 1 np.mean(cross_val_score)

Out[135]: 0.811639344262295

In [136]: 1 print(f'Heart Disease Classifier Cross-Validated Accuracy: {np.mean(cross_val_score) * 100: .2f}%')

Heart Disease Classifier Cross-Validated Accuracy: 81.16%

Area under the receiver operating characteristics curve (AUC/ROC)

- Area under curve (AUC)
- ROC curve

ROC curves are a comparison of a model's true positive (tpr) versus a model's false positive rate (fpr)

- True positive = model predicts 1 when truth is 1
- False positive = model predicts 1 when truth is 0
- True negative = models predicts 0 when truth is 0
- False negative = model predicts 0 when truth is 1

```
In [137]: 1 from sklearn.metrics import roc_curve
          2
          3 # make predictions with probabilities
          4
          5 y_probs = clf.predict_proba(x_test)
          6
          7 y_probs[:, 10], len(y_probs)
```

```
Out[137]: (array([[0.89, 0.11],
                  [0.49, 0.51],
                  [0.43, 0.57],
                  [0.84, 0.16],
                  [0.18, 0.82],
                  [0.14, 0.86],
                  [0.36, 0.64],
                  [0.95, 0.05],
                  [0.99, 0.01],
                  [0.47, 0.53]]),
          61)
```

```
In [138]: 1 y_probs_positive = y_probs[:, 1]
          2 y_probs_positive[:, 10]
```

```
Out[138]: array([0.11, 0.51, 0.57, 0.16, 0.82, 0.86, 0.64, 0.05, 0.01, 0.53])
```

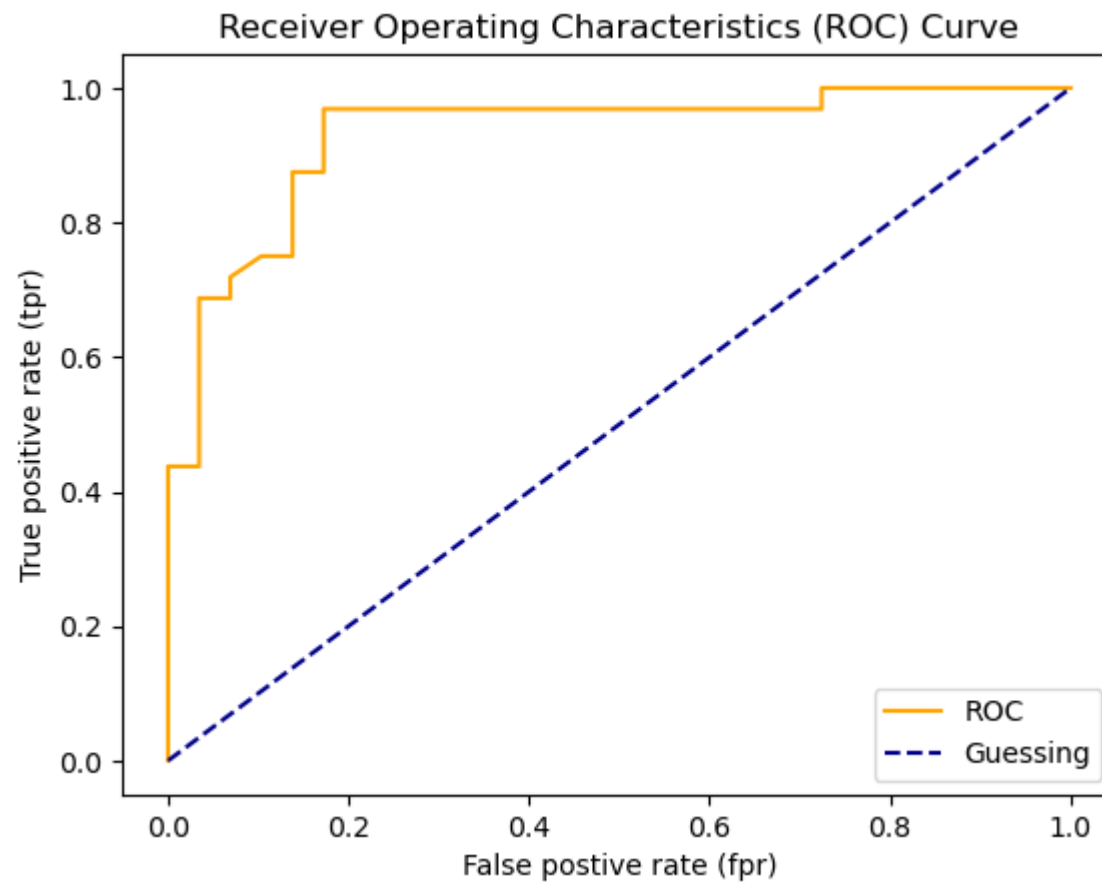
```
In [139]: 1 y_probs_positives = y_probs[:, 0]
          2 y_probs_positive[:,10]
```

```
Out[139]: array([0.11, 0.51, 0.57, 0.16, 0.82, 0.86, 0.64, 0.05, 0.01, 0.53])
```

```
In [140]: 1 # calculate fpr, tpr and thresholds
          2
          3 fpr, tpr, thresholds = roc_curve(y_test, y_probs_positive)
          4
          5 # check the false positive
          6 fpr
          7
```

```
Out[140]: array([0.          , 0.          , 0.          , 0.          , 0.          ,
                 0.03448276, 0.03448276, 0.03448276, 0.03448276, 0.06896552,
                 0.06896552, 0.10344828, 0.13793103, 0.13793103, 0.17241379,
                 0.17241379, 0.27586207, 0.4137931 , 0.48275862, 0.55172414,
                 0.65517241, 0.72413793, 0.72413793, 0.82758621, 1.          ])
```

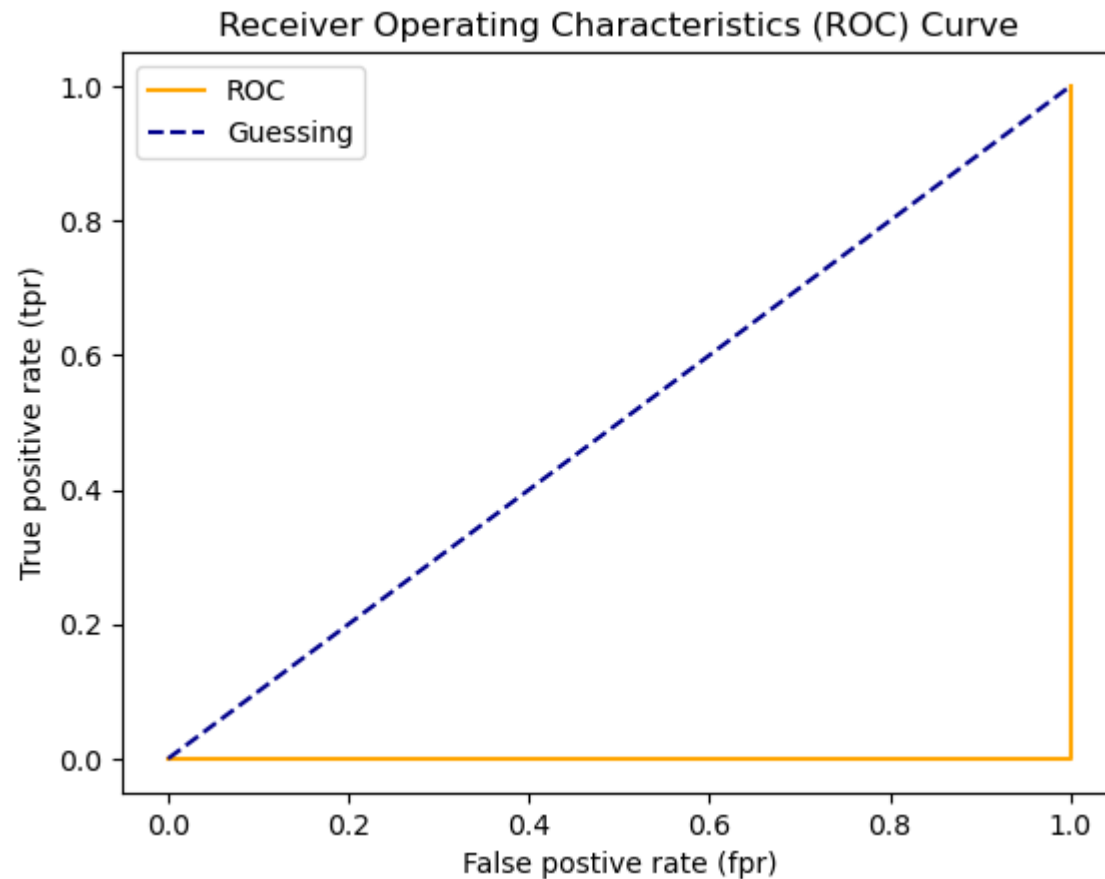
```
In [141]: 1 # create a function for plotting ROC curves
2 import matplotlib.pyplot as plt
3
4 def plot_roc_curve(fpr, tpr):
5
6     """
7     plots a Roc curve given the false positive rate (fpr)
8     and true positive rate (tpr) of a model.
9     """
10    # plot roc curves
11    plt.plot(fpr, tpr, color = 'orange', label = 'ROC')
12
13    #plot line with no predictive power (baseline)
14    plt.plot([0, 1], [0, 1], color = 'darkblue', linestyle = '--', label = 'Guessing')
15
16    #customize the plot
17    plt.xlabel('False positive rate (fpr)')
18    plt.ylabel('True positive rate (tpr)')
19    plt.title('Receiver Operating Characteristics (ROC) Curve')
20    plt.legend()
21    plt.show()
22
23 plot_roc_curve(fpr, tpr)
24
25
26
```



```
In [142]: 1 # 'Let's look at the AUC score
          2
          3 from sklearn.metrics import roc_auc_score
          4
          5 roc_auc_score(y_test, y_probs_positive)
```

Out[142]: 0.9304956896551724

```
In [143]: 1 # plot perfect ROC curve and AUC score  
2  
3 fpr, tpr, thresholds = roc_curve(y_test, y_test)  
4 plot_roc_curve(tpr, fpr)
```



```
In [144]: 1 # perfect AUC score  
2 roc_auc_score(y_test, y_test)
```

Out[144]: 1.0

Confusion matrix

A confusion matrix is a quick way to compare the label a model predicts and the actual labels it was supposed to predict.

In essence, giving you an idea of where the model is getting confused

```
In [145]: 1 from sklearn.metrics import confusion_matrix
          2
          3 y_preds = clf.predict(x_test)
          4 confusion_matrix(y_test, y_preds)
```

```
Out[145]: array([[24,  5],
                 [ 4, 28]], dtype=int64)
```

```
In [146]: 1 # visualise confusion matrix with pd.crosstab()
          2
          3 pd.crosstab(y_test,
          4             y_preds,
          5             rownames = ['Actual Labels'],
          6             colnames = ['Predicted Labels'])
```

```
Out[146]: Predicted Labels    0    1
          Actual Labels
          -----
          0    24    5
          1     4   28
```

```
In [147]: 1 24 + 5 + 4 + 28
```

```
Out[147]: 61
```

```
In [148]: 1 len(y_preds)
```

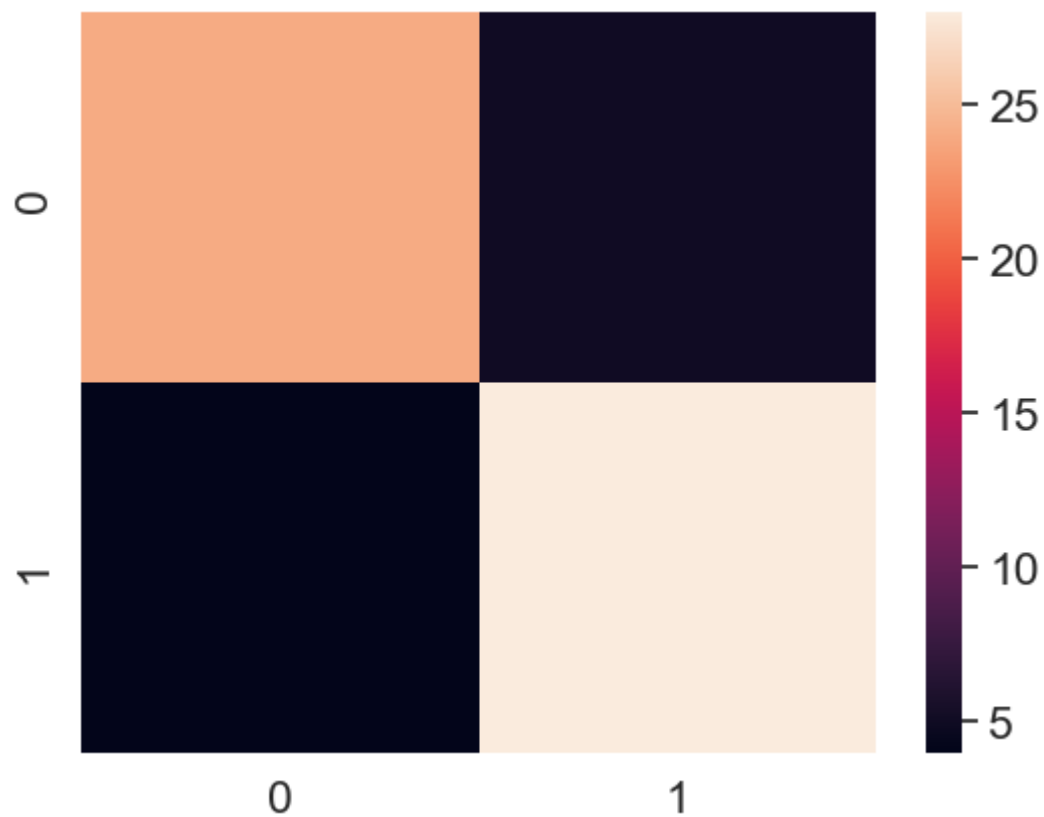
```
Out[148]: 61
```

```
In [149]: 1 len(y_test)
```

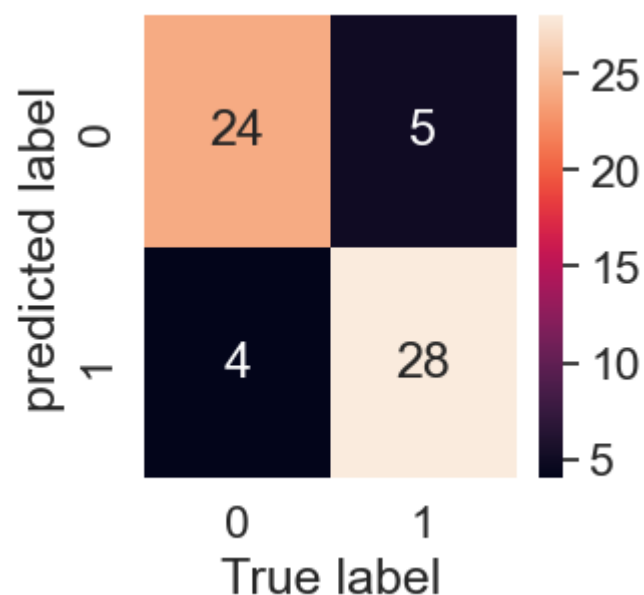
```
Out[149]: 61
```



```
In [150]: 1 # make confusion matrix more visual with seaborn's heatmap
          2 import seaborn as sns
          3
          4 # set the font scale
          5 sns.set(font_scale = 1.5 )
          6
          7 # create a confusion matrix
          8 conf_mat = confusion_matrix(y_test, y_preds)
          9
         10 # plot it using seaborn
         11 sns.heatmap(conf_mat);
```



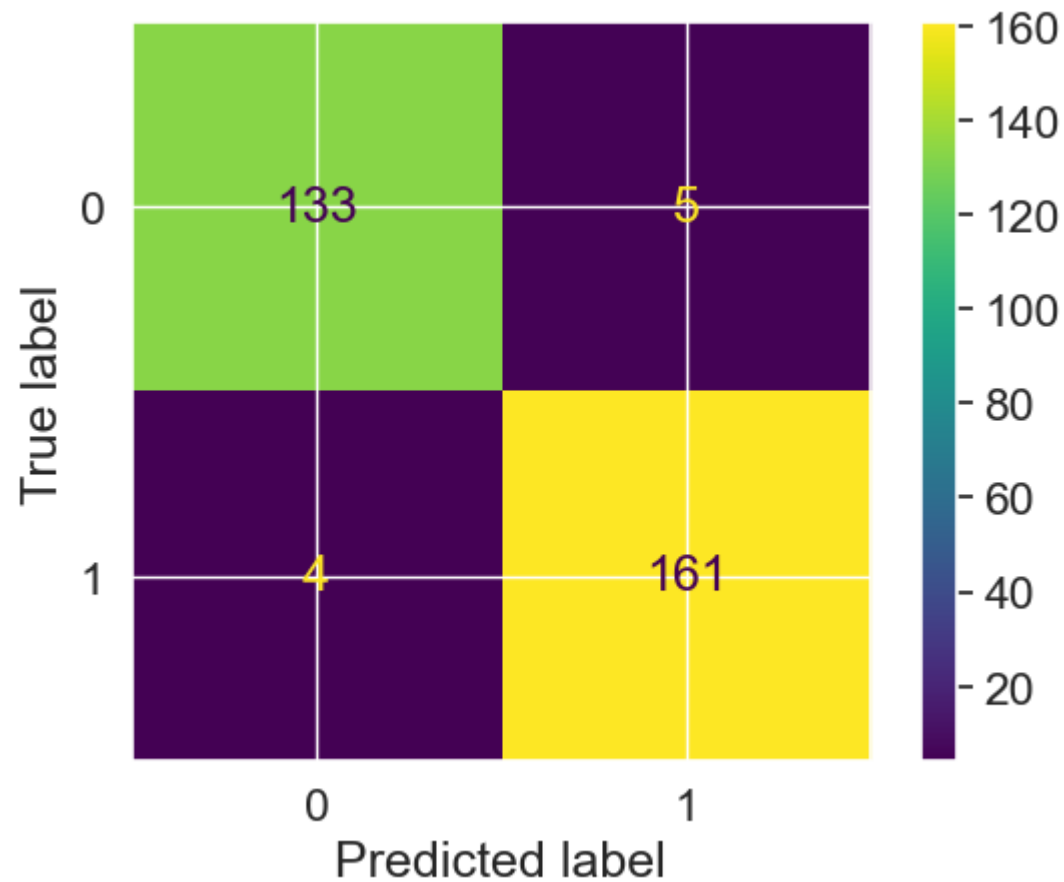
```
In [151]: 1 def plot_conf_mat(conf_mat):
2         """
3         plot a confusion matrix using seaborn's heatmap()
4         """
5         fig, ax = plt.subplots(figsize = (3, 3 ))
6         ax = sns.heatmap(conf_mat,
7                           annot = True, # Annotes the boxes with conf_mat info
8                           cbar = True)
9         plt.xlabel('True label')
10        plt.ylabel('predicted label');
11
12        plot_conf_mat(conf_mat)
```



note, if the seabon heatmap has broekn notations u can fix it with this code:

```
*bottom, top = ax.get_ylim()
*ax.set_ylim(bottom + 0.5, top -0.5)
```

```
In [152]: 1 # we can do the above using sklearn instead of seaborn(sns)
          2
          3 from sklearn.metrics import plot_confusion_matrix
          4 plot_confusion_matrix(clf, x, y);
```



Classification report

In [153]:

```
1 from sklearn.metrics import classification_report
2
3 print(classification_report(y_test, y_preds))
```

	precision	recall	f1-score	support
0	0.86	0.83	0.84	29
1	0.85	0.88	0.86	32
accuracy			0.85	61
macro avg	0.85	0.85	0.85	61
weighted avg	0.85	0.85	0.85	61

In [154]:

```
1 # Where precision and recall become valuable
2 disease_true = np.zeros(10000)
3 disease_true[0] = 1 # only one positive case
4
5 disease_preds = np.zeros(10000) # model predicts every case as 0
6
7 pd.DataFrame(classification_report(disease_true,
8                                   disease_preds,
9                                   output_dict = True))
```

Out[154]:

	0.0	1.0	accuracy	macro avg	weighted avg
precision	0.99990	0.0	0.9999	0.499950	0.99980
recall	1.00000	0.0	0.9999	0.500000	0.99990
f1-score	0.99995	0.0	0.9999	0.499975	0.99985
support	9999.00000	1.0	0.9999	10000.000000	10000.00000

To summarise classification metrics:

- **Accuracy** is a good measure to start with if all classes are balanced(e.g same amount of samples which are labelled with 0 or 1).
- **Precision** and **recall** become more important when classes are imbalanced
- if false positive predictions are worse than false negatives, aim for higher precision
- if false negative predicts are worse than false positives, aim for higher recall

- **F1-score** is a combination of precision and recall

4.2.2 Regression model evaluation metrics

1. R^2 (pronounced r-squared) or coefficient of determination
2. mean absolute error(MAE)
3. mean squared error(MSE)

R^2

what R-squared does: Compares your models predictions to the mean of the targets. Values can range from negative infinity (a very poor model) to 1. For example, if all your model does is predict the mean of the targets, it's R^2 value would be 0. And if your model perfectly predicts a range of numbers it's R^2 value would be 1.

In [155]:

```
1 boston_df.head()
```

Out[155]:

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	target
0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1.0	296.0	15.3	396.90	4.98	24.0
1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2.0	242.0	17.8	396.90	9.14	21.6
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2.0	242.0	17.8	392.83	4.03	34.7
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3.0	222.0	18.7	394.63	2.94	33.4
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3.0	222.0	18.7	396.90	5.33	36.2

```
In [156]: 1 from sklearn.ensemble import RandomForestRegressor
2
3 np.random.seed(42)
4
5 x = boston_df.drop('target', axis = 1)
6 y = boston_df['target']
7
8 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)
9
10 model = RandomForestRegressor()
11 model.fit(x_train, y_train);
```

```
In [157]: 1 model.score(x_test, y_test)
```

```
Out[157]: 0.8654448653350507
```

```
In [158]: 1 from sklearn.metrics import r2_score
2
3 # Fill an array with y_test mean
4
5 y_test_mean = np.full(len(y_test), y_test.mean())
```

```
In [159]: 1 y_test.mean()
```

```
Out[159]: 21.488235294117654
```

```
In [160]: 1 r2_score(y_test, y_test_mean)
```

```
Out[160]: 2.220446049250313e-16
```

```
In [161]: 1 r2_score(y_test, y_test)
```

```
Out[161]: 1.0
```

Mean absolute error (MAE)

MAE is the average of the absolute differences between predictions and actual values. It gives you an idea of how wrong your models predictions are.

```
In [162]: 1 # mean absolute error
          2 from sklearn.metrics import mean_absolute_error
          3
          4 y_preds = model.predict(x_test)
          5 mae = mean_absolute_error(y_test, y_preds)
          6 mae
```

Out[162]: 2.136382352941176

```
In [163]: 1 df = pd.DataFrame(data={'actual values': y_test, 'predicted values': y_preds})
          2 df['differences'] = df['predicted values'] - df['actual values']
          3 df
```

Out[163]:

	actual values	predicted values	differences
173	23.6	23.081	-0.519
274	32.4	30.574	-1.826
491	13.6	16.759	3.159
72	22.8	23.460	0.660
452	16.1	16.893	0.793
...
412	17.9	13.159	-4.741
436	9.6	12.476	2.876
411	17.2	13.612	-3.588
86	22.5	20.205	-2.295
75	21.4	23.832	2.432

102 rows × 3 columns

Mean squared error (MSE)

```
In [164]: 1 # mean squared error
          2 from sklearn.metrics import mean_squared_error
          3
          4 y_preds = model.predict(x_test)
          5 mse = mean_squared_error(y_test, y_preds)
          6 mse
```

Out[164]: 9.867437068627442

```
In [165]: 1 # Let's calculate MSE by hand
          2 squared = np.square(df['differences'])
          3 squared.mean()
```

Out[165]: 9.867437068627439

4.2.3 Finally using the scoring parameter

```
In [166]: 1 from sklearn.model_selection import cross_val_score
          2 from sklearn.ensemble import RandomForestClassifier
          3
          4 np.random.seed(42)
          5
          6 x = heart_disease.drop('target', axis = 1)
          7 y = heart_disease['target']
          8
          9 x_train, X_test, y_train, y_test = train_test_split(x, y, test_size=0.2)
          10 model = RandomForestClassifier()
```

```
In [167]: 1 np.random.seed(42)
          2 cv_acc = cross_val_score(model, x, y)
          3 cv_acc
```

Out[167]: array([0.81967213, 0.90163934, 0.83606557, 0.78333333, 0.78333333])

```
In [168]: 1 # cross-validated accuracy
          2 print(f'The cross-validated accuracy is: {np.mean(cv_acc) * 100: .2f}%')
```

The cross-validated accuracy is: 82.48%


```
In [169]: 1 # this will be the same value with the above code
          2 # here we use 'accuracy'
          3 np.random.seed(42)
          4 cv_acc = cross_val_score(model, x, y, scoring = 'accuracy')
          5 print(f'The cross-validated accuracy is: {np.mean(cv_acc) * 100: .2f}%')
```

The cross-validated accuracy is: 82.48%

```
In [170]: 1 # use of 'precision'
          2
          3 cv_precision = cross_val_score(model, x, y, scoring = 'precision')
          4 cv_precision
```

Out[170]: array([0.76315789, 0.90322581, 0.83870968, 0.79411765, 0.74358974])

```
In [171]: 1 np.mean(cv_precision )
```

Out[171]: 0.8085601538512754

```
In [172]: 1 # use of Recall
          2 cv_recall = cross_val_score(model, x, y, scoring = 'recall')
          3 np.mean(cv_recall)
```

Out[172]: 0.8424242424242424

```
In [173]: 1 # use of f1-score
          2 cv_f1 = cross_val_score(model, x, y, scoring = 'f1')
          3 np.mean(cv_f1)
```

Out[173]: 0.841476533416832

How about regression model

```
In [174]: 1 from sklearn.model_selection import cross_val_score
2 from sklearn.ensemble import RandomForestRegressor
3
4 np.random.seed(42)
5
6 x = boston_df.drop('target', axis = 1)
7 y = boston_df['target']
8
9 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)
10 rfr = RandomForestRegressor()
```

```
In [175]: 1 np.random.seed(42)
2 cv_r2 = cross_val_score(rfr, x, y, scoring = None)
3 cv_r2
```

```
Out[175]: array([0.77231143, 0.86035935, 0.74664002, 0.47632078, 0.26630379])
```

```
In [176]: 1 np.mean(cv_r2)
```

```
Out[176]: 0.6243870737930857
```

```
In [177]: 1 np.random.seed(42)
2 cv_r2 = cross_val_score(rfr, x, y, scoring = 'r2')
3 cv_r2
```

```
Out[177]: array([0.77231143, 0.86035935, 0.74664002, 0.47632078, 0.26630379])
```

```
In [178]: 1 # mean absolute error
2 cv_mae = cross_val_score(rfr, x, y, scoring = 'neg_mean_absolute_error')
3 cv_mae
```

```
Out[178]: array([-2.13045098, -2.49771287, -3.45471287, -3.81509901, -3.11813861])
```

```
In [179]: 1 # mean squared error
2
3 cv_mse = cross_val_score(rfr, x, y, scoring = 'neg_mean_squared_error')
4 cv_mse
```

```
Out[179]: array([ -7.8141513 , -12.94343325, -19.11614042, -46.28783248,
-19.48161818])
```

```
In [180]: 1 np.mean(cv_mse)
```

```
Out[180]: -21.12863512415064
```

4.3 using different evaluation metrics as scikit-learn functions.

classification evaluation functions

```
In [181]: 1 heart_disease.head(2)
```

```
Out[181]:
```

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	63	1	3	145	233	1	0	150	0	2.3	0	0	1	1
1	37	1	2	130	250	0	1	187	0	3.5	0	0	2	1

```
In [182]: 1 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
2 from sklearn.ensemble import RandomForestClassifier
3 from sklearn.model_selection import train_test_split
4
5 np.random.seed(42)
6
7 x = heart_disease.drop('target', axis = 1)
8 y = heart_disease['target']
9
10 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)
11 clf = RandomForestClassifier()
12 clf.fit(x_train, y_train)
13 y_preds = clf.predict(x_test)
14
15 # evaluate the classifier
16
17 print('Classifier metrics on the test set')
18 print(f'Accuracy: {accuracy_score(y_test, y_preds) * 100: .2f}%')
19 print(f'Precision: {precision_score(y_test, y_preds)}')
20 print(f'Recall: {recall_score(y_test, y_preds)}')
21 print(f'F1: {f1_score(y_test, y_preds)}')
```

Classifier metrics on the test set

Accuracy: 85.25%

Precision: 0.8484848484848485

Recall: 0.875

F1: 0.8615384615384615

Regression evaluation functions

```
In [183]: 1 from sklearn.metrics import r2_score, mean_absolute_error, mean_squared_error
2 from sklearn.ensemble import RandomForestRegressor
3 from sklearn.model_selection import train_test_split
4
5 np.random.seed(42)
6
7 x = boston_df.drop('target', axis = 1)
8 y = boston_df['target']
9
10 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)
11 model = RandomForestRegressor()
12 model.fit(x_train, y_train)
13 y_preds = model.predict(x_test)
14
15 # Evaluate the regression model
16 print('Regression model metrics on the test set')
17 print(f'R2: {r2_score(y_test, y_preds)* 100:.2f}%')
18 print(f'MAE: {mean_absolute_error(y_test, y_preds)}')
19 print(f'MSA: {mean_squared_error(y_test, y_preds)}')
```

Regression model metrics on the test set
R2: 86.54%
MAE: 2.136382352941176
MSA: 9.867437068627442

```
In [184]: 1 what_were_covering
```

```
Out[184]: ['0. An end-to-end scikit-learn workflow',
'1. Getting the data ready',
'2. Choose the right estimator/algorithm for our problems',
'3. Fit the models/algorithm and use it to make predictions on our data',
'4. Evaluating a model',
'5. Improve a model',
'6. Save and load a trained model',
'7. Putting it all together']
```

5. Improving a model

First predictions = baseline predictions First model = baseline model

from a data perspective:

- could we collect more data? (generally, the more data, the better)
- could we improve our data?

From a model perspective:

- is there a better model we could use?
- could we improve the current model?

Hyperparameters vs Parameters

- parameters = model find these patterns in data
- Hyperparameters = settings on a model you can adjust to (potentially) improve its ability to find patterns

Three ways to adjust hyperparameters:

1. By hand
2. Random with RandomSearchCV
3. Exhaustively with GridsearchCV

```
In [185]: 1 from sklearn.ensemble import RandomForestClassifier
          2
          3 clf = RandomForestClassifier()
```

```
In [186]: 1 # these are different hyperparameters we can adjust in RandomForestClassifier
          2 clf.get_params()
```

```
Out[186]: {'bootstrap': True,
           'ccp_alpha': 0.0,
           'class_weight': None,
           'criterion': 'gini',
           'max_depth': None,
           'max_features': 'auto',
           'max_leaf_nodes': None,
           'max_samples': None,
           'min_impurity_decrease': 0.0,
           'min_samples_leaf': 1,
           'min_samples_split': 2,
           'min_weight_fraction_leaf': 0.0,
           'n_estimators': 100,
           'n_jobs': None,
           'oob_score': False,
           'random_state': None,
           'verbose': 0,
           'warm_start': False}
```

5.1 Tuning Hyperparameters by hand

let's make 3 sets, training, validation and test.

```
In [187]: 1 clf.get_params()
```

```
Out[187]: {'bootstrap': True,  
           'ccp_alpha': 0.0,  
           'class_weight': None,  
           'criterion': 'gini',  
           'max_depth': None,  
           'max_features': 'auto',  
           'max_leaf_nodes': None,  
           'max_samples': None,  
           'min_impurity_decrease': 0.0,  
           'min_samples_leaf': 1,  
           'min_samples_split': 2,  
           'min_weight_fraction_leaf': 0.0,  
           'n_estimators': 100,  
           'n_jobs': None,  
           'oob_score': False,  
           'random_state': None,  
           'verbose': 0,  
           'warm_start': False}
```

we are going to try adjust:

- max depth
- max_feaatures
- min_samples_leaf
- min_sample_split
- n_estimators


```
In [188]: 1 def evaluate_preds(y_true, y_preds):
2         """
3         performs evaluation comparison on y_true labels vs. y_preds labels
4         on a classification model
5         """
6         accuracy = accuracy_score(y_true, y_preds)
7         precision = precision_score(y_true, y_preds)
8         recall = recall_score(y_true, y_preds)
9         f1 = f1_score(y_true, y_preds)
10        metric_dict = {'accuracy': round(accuracy, 2),
11                       'precision': round(precision, 2),
12                       'recall': round(recall, 2),
13                       'f1': round(f1, 2)}
14        print(f'Acc: {accuracy * 100: .2f}%')
15        print(f'Precision: {precision: .2f}')
16        print(f'Recalls: {recall: .2f}')
17        print(f'F1 score: {f1: .2f}')
18
19        return metric_dict
20
21
22
```

In [189]:

```
1 from sklearn.ensemble import RandomForestClassifier
2
3 np.random.seed(42)
4
5 # shuffle the data
6 heart_disease_shuffled = heart_disease.sample(frac = 1)
7
8 # split into x and y
9 x = heart_disease_shuffled.drop('target', axis = 1)
10 y = heart_disease_shuffled['target']
11
12 # split the data into train, validation and test sets
13 train_split = round(0.7 * len(heart_disease_shuffled)) # 70% of data
14 valid_split = round(train_split + 0.15 * len(heart_disease_shuffled)) # 15% of data
15 x_train, y_train = x[:train_split], y[:train_split]
16 x_valid, y_valid = x[train_split:valid_split], y[train_split:valid_split]
17 x_test, y_test = x[valid_split:], y[valid_split:]
18
19
20 clf = RandomForestClassifier()
21 clf.fit(x_train, y_train)
22
23 # make baseline predictions
24 y_preds = clf.predict(x_valid)
25
26 # evaluate the classifier on validation set
27 baseline_metrics = evaluate_preds(y_valid, y_preds)
```

Acc: 82.22%

Precision: 0.81

Recalls: 0.88

F1 score: 0.85

```
In [190]: 1 np.random.seed(42)
2
3 # create a second classifier with different hyperparameters
4 clf_2 = RandomForestClassifier(n_estimators = 100) #n_estimators by default in this version is 100 al
5 clf_2.fit(x_train, y_train)
6
7 # make predictions with different hyperparameters
8 y_preds_2 = clf_2.predict(x_valid)
9
10 # evaluate the second classifier
11 clf_2_metrics = evaluate_preds(y_valid, y_preds)
```

Acc: 82.22%

Precision: 0.81

Recalls: 0.88

F1 score: 0.85

we can adjust the parameters of the model above by hand to tune it but is a lot of work

5.2 Hyperparameter tuning with RandomizedSearchCV

In [191]:

```
1 from sklearn.model_selection import RandomizedSearchCV
2
3 grid = {'n_estimators': [10, 100, 200, 500, 1000, 1200],
4         'max_depth': [None, 5, 10, 20, 30],
5         'max_features': ['auto', 'sqrt'],
6         'min_samples_split': [2, 4, 6],
7         'min_samples_leaf': [1, 2, 4]}
8
9 np.random.seed(42)
10
11 x = heart_disease_shuffled.drop('target', axis = 1)
12 y = heart_disease_shuffled['target']
13
14 # split into train and test sets
15 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)
16
17 # Instantiate RandomForestClassifier
18
19 clf = RandomForestClassifier(n_jobs=1)
20
21 # setup RandomizedSearchCV
22 rs_clf = RandomizedSearchCV(estimator=clf, param_distributions=grid,
23                             n_iter = 10, # number of models to try
24                             cv = 5,
25                             verbose=2)
26
27 # fit the RandomizedSearchCV version of clf
28 rs_clf.fit(x_train, y_train);
```

```
Fitting 5 folds for each of 10 candidates, totalling 50 fits
[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=1200; total time= 3.8s
[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=1200; total time= 2.8s
[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=1200; total time= 2.7s
[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=1200; total time= 2.6s
[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=1200; total time= 3.0s
[CV] END max_depth=30, max_features=auto, min_samples_leaf=2, min_samples_split=4, n_estimators=100; total time= 0.2s
[CV] END max_depth=30, max_features=auto, min_samples_leaf=2, min_samples_split=4, n_estimators=100; total time= 0.2s
[CV] END max_depth=30, max_features=auto, min_samples_leaf=2, min_samples_split=4, n_estimators=100; total time= 0.3s
[CV] END max_depth=30, max_features=auto, min_samples_leaf=2, min_samples_split=4, n_estimators=100; total time= 0.3s
[CV] END max_depth=30, max_features=auto, min_samples_leaf=2, min_samples_split=4, n_estimators=100; total time= 0.2s
[CV] END max_depth=10, max_features=sqrt, min_samples_leaf=2, min_samples_split=2, n_estimators=200; total time= 0.4s
[CV] END max_depth=10, max_features=sqrt, min_samples_leaf=2, min_samples_split=2, n_estimators=200; total time= 0.5s
[CV] END max_depth=10, max_features=sqrt, min_samples_leaf=2, min_samples_split=2, n_estimators=200; total time= 0.5s
[CV] END max_depth=10, max_features=sqrt, min_samples_leaf=2, min_samples_split=2, n_estimators=200; total time= 0.5s
[CV] END max_depth=10, max_features=sqrt, min_samples_leaf=2, min_samples_split=2, n_estimators=200; total time= 0.5s
[CV] END max_depth=20, max_features=auto, min_samples_leaf=1, min_samples_split=6, n_estimators=100; total time= 0.2s
[CV] END max_depth=20, max_features=auto, min_samples_leaf=1, min_samples_split=6, n_estimators=100; total time= 0.2s
[CV] END max_depth=20, max_features=auto, min_samples_leaf=1, min_samples_split=6, n_estimators=100; total time= 0.2s
[CV] END max_depth=20, max_features=auto, min_samples_leaf=1, min_samples_split=6, n_estimators=100; total time= 0.1s
[CV] END max_depth=20, max_features=auto, min_samples_leaf=1, min_samples_split=6, n_estimators=100; total time= 0.2s
[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=1, min_samples_split=4, n_estimators=10; total time= 0.0s
```

```
[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=1, min_samples_split=4, n_estimators=10; total
time= 0.0s
[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=1, min_samples_split=4, n_estimators=10; total
time= 0.0s
[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=1, min_samples_split=4, n_estimators=10; total
time= 0.0s
[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=1, min_samples_split=4, n_estimators=10; total
time= 0.0s
[CV] END max_depth=10, max_features=auto, min_samples_leaf=2, min_samples_split=4, n_estimators=10; tota
l time= 0.0s
[CV] END max_depth=10, max_features=auto, min_samples_leaf=2, min_samples_split=4, n_estimators=10; tota
l time= 0.0s
[CV] END max_depth=10, max_features=auto, min_samples_leaf=2, min_samples_split=4, n_estimators=10; tota
l time= 0.0s
[CV] END max_depth=10, max_features=auto, min_samples_leaf=2, min_samples_split=4, n_estimators=10; tota
l time= 0.0s
[CV] END max_depth=10, max_features=auto, min_samples_leaf=2, min_samples_split=4, n_estimators=10; tota
l time= 0.0s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=500; t
otal time= 1.1s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=500; t
otal time= 1.2s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=500; t
otal time= 1.3s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=500; t
otal time= 1.1s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=500; t
otal time= 1.2s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=200; t
otal time= 0.4s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=200; t
otal time= 0.5s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=200; t
otal time= 0.5s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=200; t
otal time= 0.4s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=200; t
otal time= 0.6s
[CV] END max_depth=10, max_features=auto, min_samples_leaf=4, min_samples_split=4, n_estimators=200; tot
al time= 0.6s
[CV] END max_depth=10, max_features=auto, min_samples_leaf=4, min_samples_split=4, n_estimators=200; tot
al time= 0.5s
[CV] END max_depth=10, max_features=auto, min_samples_leaf=4, min_samples_split=4, n_estimators=200; tot
```

```
al time= 0.5s
[CV] END max_depth=10, max_features=auto, min_samples_leaf=4, min_samples_split=4, n_estimators=200; tot
al time= 0.5s
[CV] END max_depth=10, max_features=auto, min_samples_leaf=4, min_samples_split=4, n_estimators=200; tot
al time= 0.5s
[CV] END max_depth=20, max_features=sqrt, min_samples_leaf=2, min_samples_split=4, n_estimators=1000; to
tal time= 2.6s
[CV] END max_depth=20, max_features=sqrt, min_samples_leaf=2, min_samples_split=4, n_estimators=1000; to
tal time= 2.3s
[CV] END max_depth=20, max_features=sqrt, min_samples_leaf=2, min_samples_split=4, n_estimators=1000; to
tal time= 2.2s
[CV] END max_depth=20, max_features=sqrt, min_samples_leaf=2, min_samples_split=4, n_estimators=1000; to
tal time= 2.2s
[CV] END max_depth=20, max_features=sqrt, min_samples_leaf=2, min_samples_split=4, n_estimators=1000; to
tal time= 2.7s
```

In [192]: 1 rs_clf.best_params_

Out[192]: {'n_estimators': 200,
'min_samples_split': 6,
'min_samples_leaf': 2,
'max_features': 'sqrt',
'max_depth': None}

In [193]: 1 *# make predictions with the best hyperparameters*
2 rs_y_preds = rs_clf.predict(x_test)
3
4 *# evaluate the predictions*
5 rs_metrics = evaluate_preds(y_test, rs_y_preds)

Acc: 81.97%
Precision: 0.77
Recalls: 0.86
F1 score: 0.81

5.3 Hyperparameter tuning with GridSearchCV

In [194]: 1 grid

Out[194]: {'n_estimators': [10, 100, 200, 500, 1000, 1200],
'max_depth': [None, 5, 10, 20, 30],
'max_features': ['auto', 'sqrt'],
'min_samples_split': [2, 4, 6],
'min_samples_leaf': [1, 2, 4]}

```
In [195]: 1 # Let's reduce the search space for gridsearchcv/comparing with the best param
2
3 grid_2 = {'n_estimators': [100, 200, 500],
4           'max_depth': [None],
5           'max_features': ['auto', 'sqrt'],
6           'min_samples_split': [6],
7           'min_samples_leaf': [1, 2]}
```


In [196]:

```
1 from sklearn.model_selection import GridSearchCV, train_test_split
2
3 np.random.seed(42)
4
5 x =heart_disease_shuffled.drop('target', axis = 1)
6 y = heart_disease_shuffled['target']
7
8 # split into train and test sets
9 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)
10
11 # Instantiate RandomForestClassifier
12
13 clf = RandomForestClassifier(n_jobs=1)
14
15 # setup GridSearchCV
16 gs_clf = GridSearchCV(estimator=clf,
17                       param_grid=grid_2,
18                           cv =5,
19                           verbose=2)
20
21 # fit the GridSearchCV version of clf
22 gs_clf.fit(x_train, y_train);
```

```
Fitting 5 folds for each of 12 candidates, totalling 60 fits
[CV] END max_depth=None, max_features=auto, min_samples_leaf=1, min_samples_split=6, n_estimators=100; total time= 0.3s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=1, min_samples_split=6, n_estimators=100; total time= 0.4s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=1, min_samples_split=6, n_estimators=100; total time= 0.2s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=1, min_samples_split=6, n_estimators=100; total time= 0.2s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=1, min_samples_split=6, n_estimators=100; total time= 0.2s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=1, min_samples_split=6, n_estimators=200; total time= 0.8s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=1, min_samples_split=6, n_estimators=200; total time= 0.7s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=1, min_samples_split=6, n_estimators=200; total time= 0.6s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=1, min_samples_split=6, n_estimators=200; total time= 0.5s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=1, min_samples_split=6, n_estimators=200; total time= 0.7s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=1, min_samples_split=6, n_estimators=500; total time= 1.3s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=1, min_samples_split=6, n_estimators=500; total time= 1.6s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=1, min_samples_split=6, n_estimators=500; total time= 1.5s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=1, min_samples_split=6, n_estimators=500; total time= 1.5s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=1, min_samples_split=6, n_estimators=500; total time= 2.2s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=2, min_samples_split=6, n_estimators=100; total time= 0.6s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=2, min_samples_split=6, n_estimators=100; total time= 0.5s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=2, min_samples_split=6, n_estimators=100; total time= 0.4s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=2, min_samples_split=6, n_estimators=100; total time= 0.3s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=2, min_samples_split=6, n_estimators=100; total time= 0.3s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=2, min_samples_split=6, n_estimators=200; total time= 0.6s
```

```
[CV] END max_depth=None, max_features=auto, min_samples_leaf=2, min_samples_split=6, n_estimators=200; t
otal time= 0.8s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=2, min_samples_split=6, n_estimators=200; t
otal time= 0.7s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=2, min_samples_split=6, n_estimators=200; t
otal time= 0.8s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=2, min_samples_split=6, n_estimators=200; t
otal time= 0.7s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=2, min_samples_split=6, n_estimators=500; t
otal time= 1.9s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=2, min_samples_split=6, n_estimators=500; t
otal time= 1.4s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=2, min_samples_split=6, n_estimators=500; t
otal time= 1.3s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=2, min_samples_split=6, n_estimators=500; t
otal time= 1.5s
[CV] END max_depth=None, max_features=auto, min_samples_leaf=2, min_samples_split=6, n_estimators=500; t
otal time= 1.2s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=1, min_samples_split=6, n_estimators=100; t
otal time= 0.2s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=1, min_samples_split=6, n_estimators=100; t
otal time= 0.3s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=1, min_samples_split=6, n_estimators=100; t
otal time= 0.3s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=1, min_samples_split=6, n_estimators=100; t
otal time= 0.3s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=1, min_samples_split=6, n_estimators=100; t
otal time= 0.2s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=1, min_samples_split=6, n_estimators=200; t
otal time= 0.5s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=1, min_samples_split=6, n_estimators=200; t
otal time= 0.5s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=1, min_samples_split=6, n_estimators=200; t
otal time= 0.5s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=1, min_samples_split=6, n_estimators=200; t
otal time= 0.5s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=1, min_samples_split=6, n_estimators=200; t
otal time= 0.4s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=1, min_samples_split=6, n_estimators=500; t
otal time= 1.4s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=1, min_samples_split=6, n_estimators=500; t
otal time= 1.6s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=1, min_samples_split=6, n_estimators=500; t
```

```
otal time= 1.1s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=1, min_samples_split=6, n_estimators=500; t
otal time= 1.1s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=1, min_samples_split=6, n_estimators=500; t
otal time= 1.1s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=100; t
otal time= 0.2s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=100; t
otal time= 0.1s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=100; t
otal time= 0.2s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=100; t
otal time= 0.2s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=100; t
otal time= 0.1s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=200; t
otal time= 0.4s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=200; t
otal time= 0.4s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=200; t
otal time= 0.4s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=200; t
otal time= 0.4s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=200; t
otal time= 0.4s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=200; t
otal time= 0.4s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=500; t
otal time= 1.2s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=500; t
otal time= 1.2s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=500; t
otal time= 1.3s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=500; t
otal time= 1.3s
[CV] END max_depth=None, max_features=sqrt, min_samples_leaf=2, min_samples_split=6, n_estimators=500; t
otal time= 1.6s
```

```
In [197]: 1 gs_clf.best_params_
```

```
Out[197]: {'max_depth': None,  
           'max_features': 'sqrt',  
           'min_samples_leaf': 1,  
           'min_samples_split': 6,  
           'n_estimators': 200}
```

```
In [198]: 1 gs_y_preds = gs_clf.predict(x_test)  
          2  
          3 # evaluate the predictions  
          4 gs_metrics = evaluate_preds(y_test, gs_y_preds)
```

Acc: 78.69%

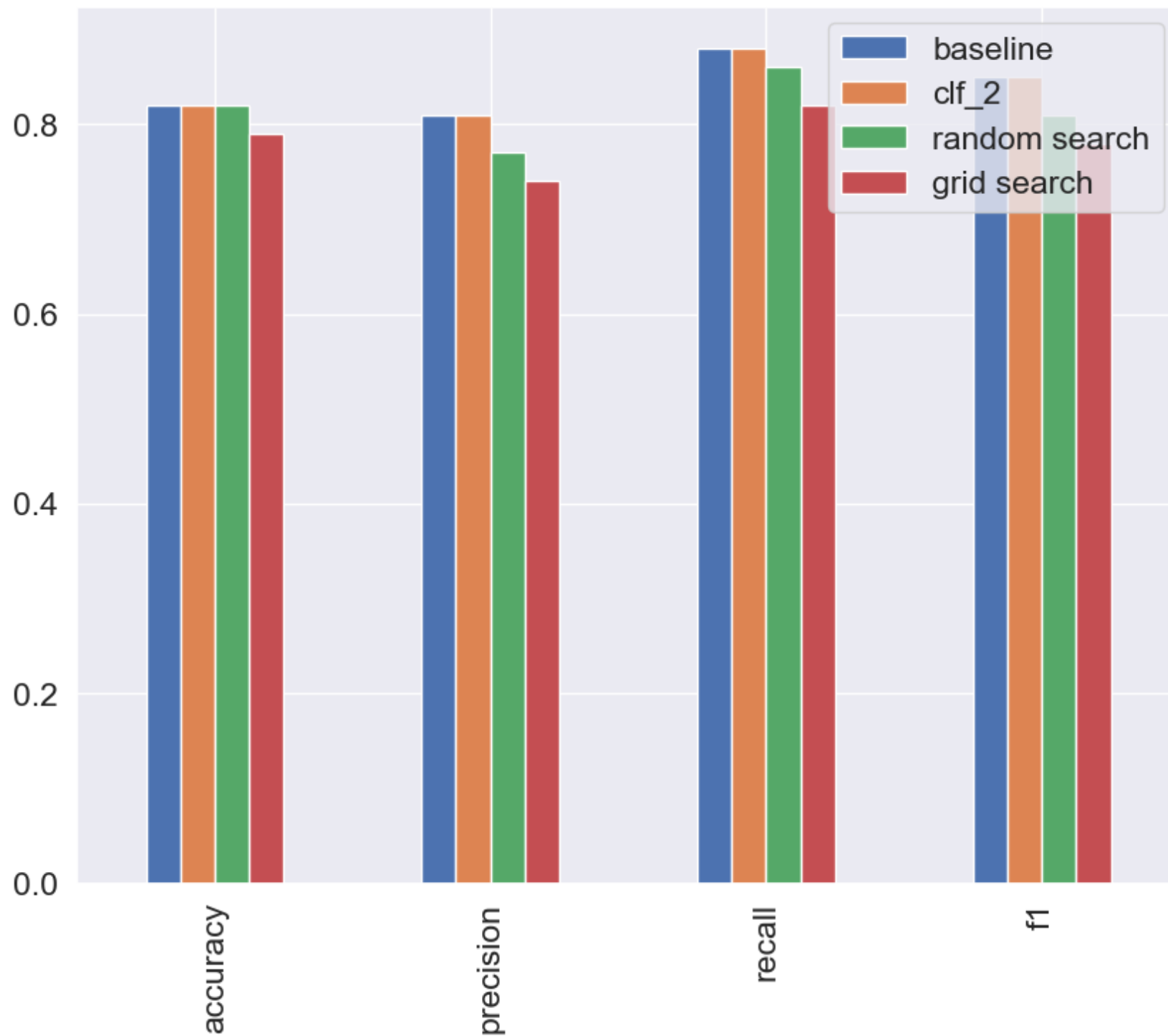
Precision: 0.74

Recalls: 0.82

F1 score: 0.78

let's compare our different models metrics

```
In [199]: 1 compare_metrics = pd.DataFrame({'baseline': baseline_metrics,  
2                                     'clf_2': clf_2_metrics,  
3                                     'random search': rs_metrics,  
4                                     'grid search': gs_metrics})  
5 compare_metrics.plot.bar(figsize = (10, 8));
```

6. Saving and loading trained machine learning models.

Two ways to save and load learning models:

1. with python's 'pickle' module
2. with the 'joblib' module

Pickle

```
In [200]: 1 import pickle
          2
          3 # save an existing model to files
          4 pickle.dump(gs_clf, open('gs_random_forest_model_1.pkl', 'wb'))
```

```
In [201]: 1 # Load a saved model
          2
          3 loaded_pickle_model = pickle.load(open('gs_random_forest_model_1.pkl', 'rb'))
```

```
In [202]: 1 # make some predictions
          2 pickle_y_preds = loaded_pickle_model.predict(x_test)
          3 evaluate_preds(y_test, pickle_y_preds)
```

Acc: 78.69%
Precision: 0.74
Recalls: 0.82
F1 score: 0.78

```
Out[202]: {'accuracy': 0.79, 'precision': 0.74, 'recall': 0.82, 'f1': 0.78}
```

Joblib

```
In [203]: 1 from joblib import dump, load
          2
          3 # save model to file
          4 dump(gs_clf, filename = 'gs_random_random_model_1.joblib')
```

```
Out[203]: ['gs_random_random_model_1.joblib']
```

```
In [204]: 1 # import a saved joblib model
          2 loaded_joblib_model = load(filename = 'gs_random_random_model_1.joblib')
```

```
In [205]: 1 # make and evaluate joblib predictions
          2 joblib_y_preds = loaded_joblib_model.predict(x_test)
          3 evaluate_preds(y_test, joblib_y_preds)
```

Acc: 78.69%
 Precision: 0.74
 Recalls: 0.82
 F1 score: 0.78

```
Out[205]: {'accuracy': 0.79, 'precision': 0.74, 'recall': 0.82, 'f1': 0.78}
```

7. Putting it all together

```
In [206]: 1 data = pd.read_csv('car-sales-extended-missing-data.csv')
          2 data
```

```
Out[206]:
```

	Make	Colour	Odometer (KM)	Doors	Price
0	Honda	White	35431.0	4.0	15323.0
1	BMW	Blue	192714.0	5.0	19943.0
2	Honda	White	84714.0	4.0	28343.0
3	Toyota	White	154365.0	4.0	13434.0
4	Nissan	Blue	181577.0	3.0	14043.0
...
995	Toyota	Black	35820.0	4.0	32042.0
996	NaN	White	155144.0	3.0	5716.0
997	Nissan	Blue	66604.0	4.0	31570.0
998	Honda	White	215883.0	4.0	4001.0
999	Toyota	Blue	248360.0	4.0	12732.0

1000 rows × 5 columns

```
In [207]: 1 data.dtypes
```

```
Out[207]: Make          object
          Colour        object
          Odometer (KM)  float64
          Doors          float64
          Price          float64
          dtype: object
```

```
In [208]: 1 data.isna().sum()
```

```
Out[208]: Make          49
          Colour        50
          Odometer (KM)  50
          Doors          50
          Price          50
          dtype: int64
```

Steps we want to do (all in one cell):

1. fill missing data
2. convert the data to numbers
3. build a model on the data


```
In [1]: 1 # getting started
2 import pandas as pd
3 from sklearn.compose import ColumnTransformer
4 from sklearn.pipeline import Pipeline
5 from sklearn.impute import SimpleImputer
6 from sklearn.preprocessing import OneHotEncoder
7
8 # modelling
9 from sklearn.ensemble import RandomForestRegressor
10 from sklearn.model_selection import train_test_split, GridSearchCV
11
12 # set random seed
13 import numpy as np
14 np.random.seed(42)
15
16 # import data and drop rows with missing labels
17 data = pd.read_csv('car-sales-extended-missing-data.csv')
18 data.dropna(subset=['Price'], inplace=True)
19
20 # define different features and transform pipeline
21 categorical_features = ['Make', 'Colour']
22 categorical_transformer = Pipeline(steps=[
23     ('imputer', SimpleImputer(strategy='constant', fill_value='missing')),
24     ('onehot', OneHotEncoder(handle_unknown='ignore'))])
25
26 door_feature = ['Doors']
27 door_transformer = Pipeline(steps=[
28     ('imputer', SimpleImputer(strategy='constant', fill_value=4))
29 ])
30
31 numeric_features = ['Odometer (KM)']
32 numeric_transformer = Pipeline(steps=[
33     ('imputer', SimpleImputer(strategy='mean'))
34 ])
35
36 # setup preprocessing steps (fill missing values, then convert to numbers)
37 preprocessor = ColumnTransformer(
38     transformers = [
39         ('cat', categorical_transformer, categorical_features),
40         ('door', door_transformer, door_feature),
41         ('num', numeric_transformer, numeric_features)
42     ])
43
```

```
44 # creating a preprocessing and modeling pipeline
45 model = Pipeline(steps=[('preprocessor', preprocessor),
46                           ('model', RandomForestRegressor())])
47
48 # split data
49 x = data.drop('Price', axis = 1)
50 y = data['Price']
51 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)
52
53 # fit and score the model
54 model.fit(x_train, y_train)
55 model.score(x_test, y_test)
56
```

Out[1]: 0.22188417408787875

it's also possible to use GridSearchCV or RandomizedSearchCV with our Pipeline

```

In [4]: 1 # use GridSearchCV with our regression Pipeline
        2 from sklearn.model_selection import GridSearchCV
        3
        4 pipe_grid = {
        5     'preprocessor__num__imputer__strategy': ['mean', 'median'],
        6     'model__n_estimators': [100, 1000],
        7     'model__max_depth': [None, 5],
        8     'model__max_features': ['auto'],
        9     'model__min_samples_split': [2, 4]
       10 }
       11
       12 gs_model = GridSearchCV(model, pipe_grid, cv=5, verbose=2)
       13 gs_model.fit(x_train, y_train)

```

Fitting 5 folds for each of 16 candidates, totalling 80 fits

```

[CV] END model__max_depth=None, model__max_features=auto, model__min_samples_split=2, model__n_estimators=100, preprocessor__num__imputer__strategy=mean; total time= 0.4s
[CV] END model__max_depth=None, model__max_features=auto, model__min_samples_split=2, model__n_estimators=100, preprocessor__num__imputer__strategy=mean; total time= 0.3s
[CV] END model__max_depth=None, model__max_features=auto, model__min_samples_split=2, model__n_estimators=100, preprocessor__num__imputer__strategy=mean; total time= 0.3s
[CV] END model__max_depth=None, model__max_features=auto, model__min_samples_split=2, model__n_estimators=100, preprocessor__num__imputer__strategy=mean; total time= 0.3s
[CV] END model__max_depth=None, model__max_features=auto, model__min_samples_split=2, model__n_estimators=100, preprocessor__num__imputer__strategy=mean; total time= 0.4s
[CV] END model__max_depth=None, model__max_features=auto, model__min_samples_split=2, model__n_estimators=100, preprocessor__num__imputer__strategy=median; total time= 0.4s
[CV] END model__max_depth=None, model__max_features=auto, model__min_samples_split=2, model__n_estimators=100, preprocessor__num__imputer__strategy=median; total time= 0.4s
[CV] END model__max_depth=None, model__max_features=auto, model__min_samples_split=2, model__n_estimators=100, preprocessor__num__imputer__strategy=median; total time= 0.5s
[CV] END model__max_depth=None, model__max_features=auto, model__min_samples_split=2, model__n_estimators=100, preprocessor__num__imputer__strategy=median; total time= 0.4s
[CV] END model__max_depth=None, model__max_features=auto, model__min_samples_split=2, model__n_estimators=100, preprocessor__num__imputer__strategy=median; total time= 0.4s

```

```

In [5]: 1 gs_model.score(x_test, y_test)

```

Out[5]: 0.3339554263158365

```

In [ ]: 1

```

